# A Survey:-Optimize Visual Image Search for Semantic web

Jaimin Shroff  , Prof. Parth D. Shah  , Prof. Bhaumik Nagar

*Charotar University of Science and Technology,(CHARUSAT), changa, Anand, Gujarat-388421*

**Abstract:-In today's web image search engines find more irrelevancies in the search result. By adding semantic meaning to the document this relevancy can be eliminated. SIEVE image search algorithm combine the text based and content based method and shows the result. Also "IN-Picture" search algorithm mixing the images higher level and lower level contents. In this paper it shows some image searching framework like "SAFE" describe how image are searched using its attributes. Also describes some sematic web technology, which helps in image search and shows how detailed indexing system can use SPRQL query and ontology of an image to build semantic web based framework.**

**Keyword: - RDF, SIEVE, Semantic, SPRQL, Ontology, OWL**

## 1. INTRODUCTION:-

Web 1.0 is a system of interlinked, hypertext documents accessed via the Internet. With a Web browser, a user views Web pages that may contain text, images, and other multimedia and navigates between them using hyperlinks.Web2.0 is a perceived second generation of web-based communities and hosted services such as social-networking sites, wikis which facilitate collaboration and sharing between users. Semantic Web (Web 3.0) is not a separate Web but an extension of the Web 2.0 in which information given with well-defined meaning. It is used to purport the useful information from the web. Semantics mean adding meaning of data to be discovered by computers. The Semantic Web is a vision: the idea of having data on the Web defined and linked in such a way that it can be used by machines not just for display purposes, but for automation, integration and reuse of data across various applications. Semantic Web technologies enable people to create data stores on the Web, build vocabularies, and write rules for handling data.

## 2. CURRENT IMAGE SEARCH ENGINES:-

Although Internet has contributed a lot for human society, the explosive growth of multimedia data transmission has generated a critical need for efficient, high-capacity image databases, as well as powerful search engines to retrieve image desired from them. At present, main search engines have developed a series of feasible solutions for text search as text is a kind of structural information, which also have already been taken into commercial use phrase. Google PageRank algorithm, has reached satisfactory result on text information search area, helping users find useful information in relatively short time, but with the pluralism of internet media, unique text search service cannot meet customers' current requirement. As one of the most important media for presentation, images have become significant as text information recently, but the performance of images search engine is still less than satisfactory.

An image search is a systematic process that includes browsing, searching and retrieving images from a large database of digital images.

### 2.1Different Frameworks of Image search:-

*2.1.1 Framework for Picture Extraction on Search Engine*: it finds the image of the base of the higher level and lower level concept features of an image.[6] Also using some different approach of searching techniques like text base search content based search, context based search etc.

*2.2.2 SaFe: A General Framework for Integrated Spatial and Feature Image Search*: it finds the image based on its attributes of regions. Features of attributes includes color, shape, texture etc. [7] for the images.it demonstrate the spatial features for an images.
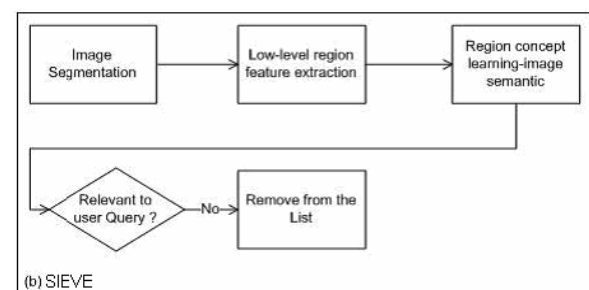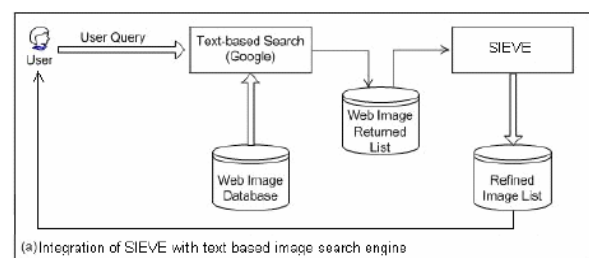
### 2.2SIEVE Algorithm for effective image search:-



Fig:- 2.1Block diagram of image retrieval system and sieve algorithm

Fig.2.1 Show the block diagrams of the proposed retrieval system and the SIEVE module. The input to the system is a keyword query submitted by the user, the system first comes out a ranked list of retrieved images given by the text based image search engine. The list of retrieved images by text based image search engine is then input to SIEVE for further analysis. For each image in the list, SIEVE first segments it into different regions. Then, color and texture features of each region are extracted using the methods presented in [11], [12] and [13], [14] respectively. The region color feature is the dominant color in HSV space and the region texture feature is the Gabor feature obtained using a novel padding algorithm.

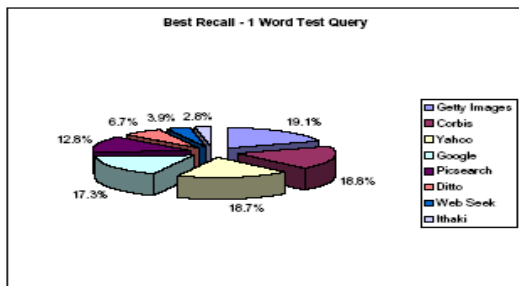## 2.3 Comparative study of Image Search engines:-
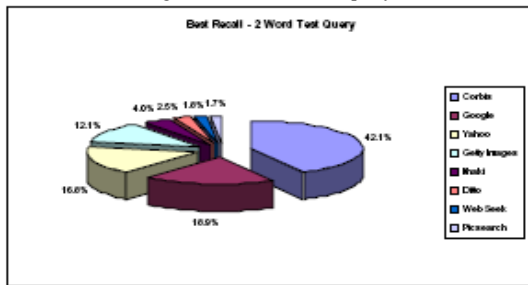


Fig. 2.3.1 one word test query



Fig 2.3.2 Two word test query

Above two figures describe the graph of the comparative study of the different search engine [14].It shows that how different search engine give more accurate result when the query is sort. But when the words in the query are incremented some of the search engine's accuracy is decremented.it shows that when query is too long the meaning is divided and the accuracy is decremented.

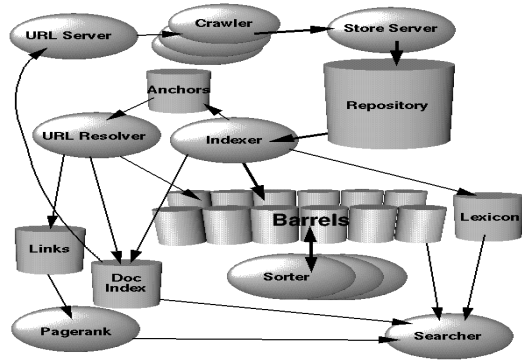## 2.4 Search Engine Architecture Overview:-



Fig.2.4.1 Search engine Indexing System

In Search engines, the web crawling is done by several distributed crawlers. There is a URLserver that sends lists of URLs to be fetched to the crawlers. The web pages that are fetched are then sent to the storeserver. The storeserver then compresses and stores the web pages into a repository. Every web page has an associated ID number called a docID which is assigned whenever a new URL is parsed out of a web page. The indexing function is performed by the indexer and the sorter. The indexer performs a number of functions. It reads the repository, uncompressed the documents, and parses them. Each document is converted into a set of word occurrences called hits. The hits record the word, position in document, an approximation of font size, and capitalization. The indexer distributes these hits into a set of barrels creating a partially sorted forward index. The indexer performs another important function. It parses out all the links in every web page and stores important information about them in an anchors file. This file contains enough information to determine where each link points from and to, and the text of the link.

The URLresolver reads the anchors file and converts relative URLs into absolute URLs and in turn into docIDs. It puts the anchor text into the forward index, associated with the docID that the anchor points to. It also generates a database of links which are pairs of docIDs. The links database is used to compute Page Ranks for all the documents.

The sorter takes the barrels, which are sorted by docID and resorts them by wordID to generate the inverted index. This is done in place so that little temporary space is needed for this operation. The sorter also produces a list of wordIDs and offsets into the inverted index. A program called Dump Lexicon takes this list together with the lexicon produced by the indexer and generates a new lexicon to be used by the searcher. The searcher is run by a web server and uses the lexicon built by Dump Lexicon together with the inverted index and the page Ranks to answer queries.

2.4.1 Repository:-

The repository contains the full HTML of every web page. Each page is compressed using zlib.

In the repository, the documents are stored one after the other and are prefixed by docID, length, and URL.

2.4.2 Document Index:-

The document index keeps information about each document. The information stored in each entry includes the current document status, a pointer into the repository, a document checksum, and various statistics. If the document has been crawled, it also contains a pointer into a variable width file called doc info which contains its URL and title. Otherwise the pointer points into the URLlist which contains just the URL.

2.4.3 Hit Lists:-

A hit list corresponds to a list of occurrences of a particular word in a particular document including position, font, and capitalization information. Hit lists account for most of the space used in both the forward and the inverted indices. Because of this, it is important to represent them as efficiently as possible.

2.4.4Forward Index:-

The forward index is actually already partially sorted. It is stored in a number of barrels. Each barrel holds a range of word ID's. If a document contains words that fall into a particular barrel, the docID is recorded into the barrel, followed by a list of word ID's with hitlists which correspond to those words. This scheme requires slightly more storage because of duplicated docIDs but the difference is very small for a reasonable number of buckets and saves considerable time and coding complexity in the final indexing phase done by the sorter. Furthermore, instead of storing actual word ID, we store each wordID as a relative difference from the minimum wordID that falls into the barrel the wordID is in.

2.4.5Inverted Index:-

The inverted index consists of the same barrels as the forward index, except that they have been processed by the sorter. For every valid wordID, the lexicon contains a pointer into the barrel that wordID falls into. It points to a doclist of docID's together with their corresponding hit lists. This doclist represents all the occurrences of that word in all documents.

Algorithm of the Search engine:-

1.      Parse the query.
2.      Convert words into wordIDs.
3.      Seek to the start of the doclist in the short barrel for every word.
4.      Scan through the doclists until there is a document that matches all the search terms.
5.      Compute the rank of that document for the query.
6.      If we are in the short barrels and at the end of any doclist, seek to the start of the doclist in the full barrel for every word and go to step 4.
7.      If we are not at the end of any doclist go to step 4.
Sort the documents that have matched by rank and return the top k.

2.4.6Anchor Text:-

The text of links is treated in a special way in our search engine. Most search engines associate the text of a link with the page that the link is on. Anchors often provide more accurate descriptions of web pages than the pages themselves. Anchors may exist for documents which cannot be indexed by a text-based search engine, such as images, programs, and databases. This makes it possible to return web pages which have not actually been crawled.

### 3. DIFFERENT SEMANTIC WEB TECHNOLOGIES:-

The Semantic Web extends the Web through the use of standards, markup languages and related processing tools. Semantic Web technologies enable people to create data stores on the Web, build vocabularies, and write rules for handling data. Linked data are empowered by technologies such as RDF, SPARQL, and OWL.

*3.1RDF (Resource description Framework):-*

The RDF (Resource Description Framework) is a language for describing information and resources on the web. Putting information into RDF files, makes it possible for computer programs to search, discover, pick up, collect, analyze and process information from the web. It generates a model of all the documents is called schemas [iii].

Web Resources:-A web resource is simply any identifiable information on the web. The resource itself is conceptual, while its representation is actual. When a web resource is requested, an appropriate representation of its current state is provided.

*3.2Ontology:-*

Ontology formally represents knowledge as a set of concepts within a domain, and the relationships between those concepts. It can be used to reason about the entities within that domain and may be used to describe the domain. A conceptualization refers to an abstract model of some phenomenon in the world by identifying the relevant concept of that phenomenon. Explicit means that the types of concepts used and the constraints on their use are explicitly defined.

*3.3OWL:-*

The OWL Web Ontology Language is designed for use by applications that need to process the content of information instead of just presenting information to humans. OWL facilitates greater machine interpretability of Web content than that supported by XML, RDF, and RDF Schema (RDF-S) by providing additional vocabulary along with a formal semantics [iv].
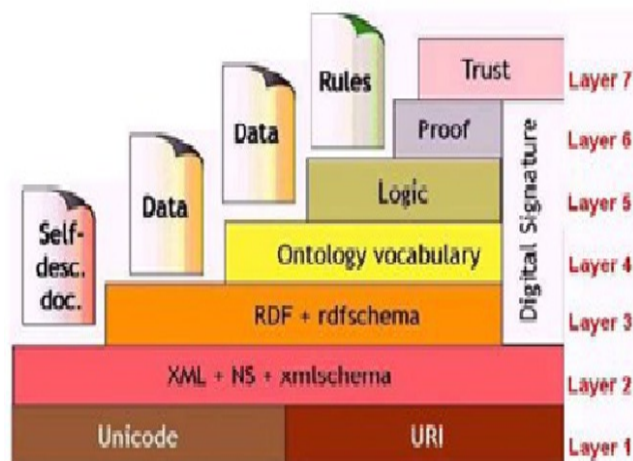
### 4. STRUCTURE OF SEMANTIC WEB:-



Fig. 4.1 Sematic Web Architecture

The Unicode and Uniform Resource Identifier (URI) layers make sure that international characters sets are used and provide means for identifying the objects in the Semantic Web. The XML layer with namespace and schema definitions make sure the Semantic Web definitions can integrate with the other XML based standards. XML provides a surface syntax for structured documents, but imposes no semantic constraints on the meaning of these documents.

The Resource Description Framework (RDF) presents a simple model that can be used to represent any kind of data. This data model consists of nodes connected by labeled arcs, where the nodes represent Web resources and the arcs represent properties of these resources [v].

### 5. CHALLENGES IN EFFECTIVE IMAGE SEARCH:-

- Vocabulary: what kinds of image features should be used? How to map them to words? The most generally utilized method is clustering. Some researchers also adopted a hierarchical clustering method to generate a vocabulary tree. But it is clear that we need to develop some kinds of visual language models to solve the problem.

- Long query: The reason why text search engine is effective is because text queries usually only contain a few words. So, the query document matching can be conducted efficiently by inverted index. Although images can be represented by "bag of features," the retrieval problem is still very different from text retrieval because query-by example is actually equivalent to using a whole document as a query. So, the search is more like document to document matching.

- Content quality: Web search engine is effective because it can use link analysis to obtain quality and importance measurement for Web pages. For images, it is hard to obtain similar kind of measurement because the links are typically not directly associated with images. Without PageRank for images, we won't be able to take advantage of many top-k search techniques typically used in web search, and it also will lead to the lack of efficient cache of index.

- Relevance ranking: The similarity measure between two images is quite different from text. How image words are weighted in computing the relevance. And how to deal with "word proximity" in images [16].

### 6. PROPOSED FRAMEWORK FOR EFFECTIVE IMAGE SEARCH:-

To find the best solution related to our query of our image search techniques is not enough? We have to use the sematic web technology it provide more detailed meaning of data. Semantic web provides more structured data than the current web. When lexicon finds the different words of the query in the repository we provide an ontology based SPRQL query between lexicon and repository. That will help machine to find the classes and subclass of the related word. Making the ontology of an image means to classify an image in detail using its contents.

### REFERENCES:-

I. Google image search: http://images.google.com
II. Yahoo image search: http://images.yahoo.com
III. http://www.w3.org/RDF/
IV. http://www.semanticalley.com/category/web-semantic/
V. http://www.webcentralstation.ca/2011/02/08/an-intro-to-the-semantic-web-why-you-need-to- know-about-it-sooner-than-later/
VI. http://www.swoogle.umbc.edu/

*Number of Technical Paper Referred:-*

[1] "SIEVE—Search Images Effectively through Visual Elimination" By Ying Liu, Dengsheng Zhang and Guojun LuGippsland School of Info Tech, Monash University, Churchill, Victoria, 3842

[2] "PageRank for Product Image Search" by Yushi Jing, Shumeet Baluja College Of Computing, Georgia Institute of Technology, Atlanta GA 2Google, Inc. 1600 Amphitheater Parkway, Mountain View, CA

[3] "Speeded up robust features "by H. Bay, T. Tuytelaars, and L. V. Gool. Surf.

[4] "Shape matching and object recognition using shape contexts "by S. Belongie, J. Malik, and J. Puzicha.

[5] "An image and video search engine for the World-Wide Web". By Smith, J. R. and Chang, S.-F. 1997. Jose, CA), 84 95.

[6] "A Framework for Picture Extraction on Search Engine Improved and Meaningful Result" by Anamika Sharma, Sarita Sharma Computer Science, MDU Rohtak, DAVIM Faridabad, Haryana, India

[7] "SaFe: A General Framework for Integrated Spatial and Feature Image Search" by john R. Smith Shih-Fu Chang Dept. of Electrical Engineering Dept. of Electrical Engineering Columbia University Columbia University New York,

[8] "A flexible approach for managing digital images on the semantic web" by Halaschek-Wiener, C., Schain, A., Golbeck, J., Grove, M., Parsia, B., Hendler, J.

[9] "Hierarchical Clustering of WWW Image Search Results using Visual, Textual and Link Information" by D. Cai, X. He, Z. Li, W. -Y. Ma and J. -R. Wen.

[11] "Boosting Web Image Search by Co- Ranking" by J. He, C. Zhang, N. Zhao and H. Tong.

[12] "Region-Based Image Retrieval with High-Level Semantic Color Names" Y. Liu, D. S. Zhang, G. Lu, and W.-Y. Ma.

[13] "Region-based Image Retrieval with Perceptual Colors "by Y. Liu, D. S. Zhang, G. Lu and W. -Y. Ma.

[14] "Study on Texture Feature Extraction in Region-Based Image Retrieval System" by Y. Liu, D. S. Zhang, G. Lu and W. -Y. Ma.

[15] "Comparative evaluation of web image search engines for multimedia application "By Keon Stevenson and Clement Leung, School of Computer Science and Mathematics Victoria University,

[16] "Searching One Billion Web Images by Content: Challenges and Opportunities" by Zhiwei Li, Xing Xie, Lei Zhang, and Wei-Ying Ma Web Search and Data Mining Group Microsoft Research Asia