

## Smart Crawling System- A Review

Parmar Krina

Department of Computer Engineering  
MET's Institute of Engineering  
Nashik, India

Misal Kiran

Department of Computer Engineering  
MET's Institute of Engineering  
Nashik, India

Nemade Neha

Department of Computer Engineering  
MET's Institute of Engineering  
Nashik, India

Kulkarni Prajka

Department of Computer Engineering  
MET's Institute of Engineering  
Nashik, India

---

**Abstract:** *The sketch-based image retrieval is a system in which we input an image. Usually there is a huge appearance gap in sketches and photo-realistic images, so to bridge the gap we design a framework that is histogram of line relationship (HLR) which is line segment-based descriptor and a noise reduction algorithm called as Object boundary selection. Another focus is on Smart crawler (input in form of text), in which the deep web grows very rapidly, therefore high efficiency and wide coverage is a major issue. To overcome such challenging issue, a two stage Smart Crawler can be used. We propose a system, in which we are going to combine two systems namely Sketch-based image retrieval and Smart Crawler. The SBIR not only works on the colors but also focuses on the shapes and boundaries. This will provide an efficient image by removing the noisy edges. In the initial stage of Smart Crawler, to achieve accurate results site-based searching is done and it ranks the websites to give priorities to the relevant ones. In second stage, smart crawler achieves site searching by digging up useful links in link-ranking. Hence we obtain the useful information by avoiding irrelevant links.*

**Keywords:** *HLR, Object boundary Selection, Deep web, Link-ranking.*

---

### I. INTRODUCTION

Instead of writing, using sketches and drawing the images are better for easy understanding and communication. However, the oldest form of writing was logo-graphic method. Shapes give the rough shape of an object and also provide conceptual representation for easy communication. Objects can be easily recognized from other sketches which provide the tools for people for communicating in different languages. Instead of writing, drawing sketches are more natural and informative. One of the most valuable tools, called smart crawling system which provides the sketch-based image retrieval, which is supplementary for keywords-based search. The ever growing applications of sketch-based retrieval can be seen in touch-screen devices like smart devices, tablets. Though sketches are considered good way of expressing people's thoughts, there is a difference between the user's sketches and the photo-realistic images. The structured of an object is mainly focused by the people and they only draw the semantic contour boundary, whereas photo-realistic images contain the color, texture and details of the shape of an object and many more detailed information. Thus comparing the user-sketches and the photo-realistic images are difficult. The alternative way for achieving this is to use edge extraction technique. The edge extraction can be applied on a photo-realistic image before matching. After edge extraction, the photo-realistic images are represented by strong edges; through which comparison can be simplified. The extracted edges are set of lines. To capture line level features, descriptors must be capable. Line-based descriptors give the flexibility for edge selection or removal. This is done by setting the corresponding parts of the feature vector at a certain value, which is critical for boundary selection. Some of the existing descriptors for sketched-based image retrieval are SHoG, GF-HOG. These are designed for capturing the pixel-level features from the patches of images.

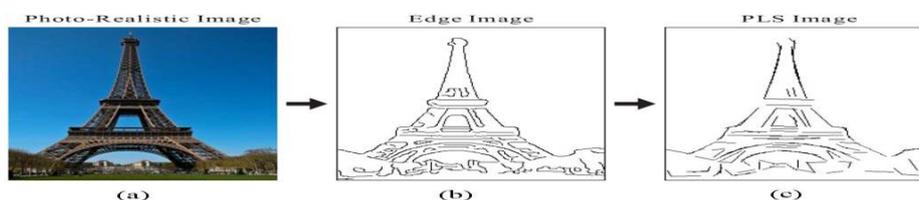


Fig.1. Image preprocessing. [11]

(For photo-realistic image, by applying Canny edge detector strong edges can be extracted, and these extracted edges are then approximated into a set of line segments.) Along with images, sketches, web also contains large volume of web resources. But due to dynamic nature, to achieve wide coverage and high efficiency is a challenging task. However, sometimes the hidden web pages are highly visited by some highly relevant links. The deep web refers to contents behind the searchable web interfaces that are not indexed by search engines. Thus, deep web refers to the contents that lie behind the searchable interfaces that are not indexed by probing engines.

Predictions on extrapolations are done from a study at University of California, Berkeley, it is estimated that the deep web contains 91,850 tb and the surface web is only about 167 tb in the year 2003. The recent studies estimated that 1.9 zb were reached and 0.3 zb were consumed ecumenical in the year 2007. An IDC report estimates that the total of all digital data created, replicated, and consumed will reach 6 zb in the year 2014. A consequential portion is large amount of data that is estimated to be stored as structured or relational data in databases. Deep web makes up about 96% of all the content on the Internet, which is 500 to 550 times more immense than the surface web. Thus these data contains an astronomical amount of valuable information and entities such as Info mine, Crusty, Books may be fascinated with building an index of the deep web sources in a given domain. There is a need for an efficient crawling system that will be able to accurately and quickly explore the deep web databases. A new challenge is to locate the deep web databases, because they are not registered with any search engines, are usually sparsely distributed, and keep constantly changing. To address this problem, previous work proposed two types of crawlers namely, generic crawlers and focused crawlers [12].

## II. RELATED WORK

### A. Sketch based Image Retrieval:

Tianjia shao et al proposed, that the increasing in complexity of geometric models, leads to new challenges for sketch-based shape retrieval systems. The use of vectorized contours for 2D shape representation was introduced. We have also developed a robust sampling-based shape matching algorithm. But there are several limitations of this work that are relatively easy to address. They only use contour images from seven viewpoints to compare with query sketch [1]. Huda A. Abdulbaqi et al discusses about the CBIR in reaching specifically to the specific ways depending on low level (shape, color, texture) and high level (including semantic). The SBIR is discussed in base of indexing, feature extraction, matching and geometrical element (rotation, scaling, transformation). For large scale sketch, TENSOR based image descriptor developed to be superior comparing with MPEG-7 EHD. But MPEG-7 EHD. SBIR system that used SHoG descriptor as a better in benchmark [2]. Neetesh Prajapati and G.S.Prajapati considered the two main features in the report. The recovery procedure has to be very unusual and interactive. The paper presents the dissimilar techniques used to execute, implement plan & examination a sketch based image retrieval system. HOG is more effective than the EHD [4]. Ms. Asmita A. Desai et al describes on the individual sketch method and try to integrate this method for reducing the individual methods drawback. In this technique EHD and HOG method work parallel and normalized result is displayed. Integrated system gives the better accuracy [6]. Prachi A. Gaidhani and S. B. Bagal state the concept of content-based and sketch based image retrieval in this paper. Future research directions have been suggested and open research issues were identified in this paper [7].

### B. Smart Crawler:

Radhika Bairagade et al proposed an effective harvesting framework for deep-web interfaces, namely Smart Crawler. Author had proven that approach achieves both wide coverage for deep web interfaces and maintains highly efficient crawling.[3]. Sneha A. Ghumatkar et al contains different kind of general searching technique and Meta search engine strategy and by using it an efficient way of searching most

relevant data from hidden web. It combined multiple search engine and two stage crawler for harvesting most relevant site.[5]. Nimisha Jain et al discuss about the many previous challenges such as efficiency, end-to-end delay, quality of link, failure to find the deep websites as they are unregistered with any crawler, scattered and dynamic. This approach accomplishes wide spread coverage of deep web and implements proficient crawling technique. They add more dataset to expand the deep web dataset to crawl more number of websites giving maximum possible search results [8]. Archana Pondkule et al proposed a Smart Crawler that performs site-based locating by reverse searching method is known deep web sites for center pages. They can effectively find more data sources for sparse domains [9]. Ms. Asmita D Rathod had proven that approach achieves both wide coverage for deep web interfaces and maintains highly efficient crawling. Experimental results on a representative set of domains shows the effectiveness of the proposed two-stage crawler, which achieves higher harvest rates than other crawlers.[10]

### III. PROPOSED SYSTEM

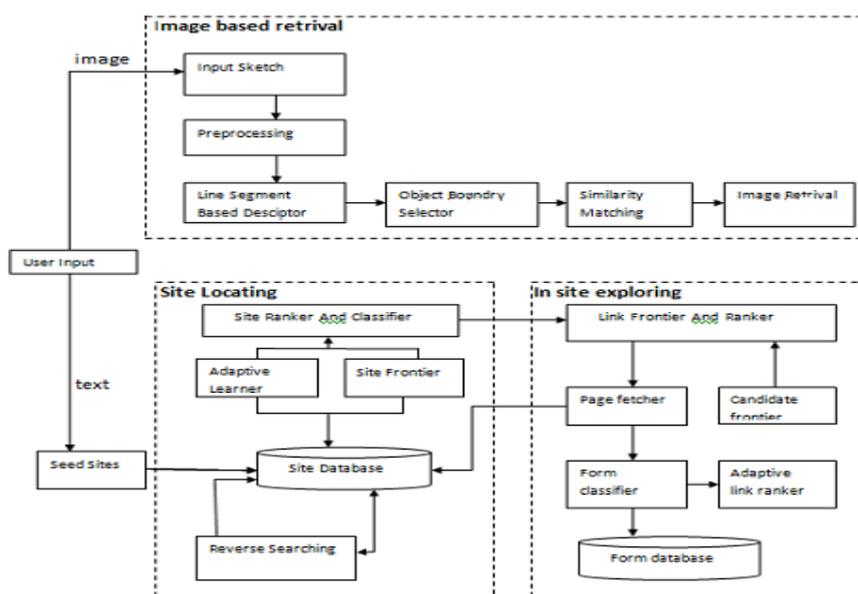


Fig 2. Architecture of smart crawling system

Our mainly important task is to overpass the information gap between the drawing and the image, which is assist by own pre-processing alteration process. The repetition of the consumption process is probable, by the existing results looking again, thus increasing the precision. The system building blocks consist a pre-processing subsystem, which remove the problems caused by the plurality of metaphors. Using the attribute vector generating subsystem our image can be represented by numbers considering a given property. The database management subsystem provides a medium between the database and the agenda. Bottom on the feature vectors and the model image the retrieval subsystem provides the response list for the user using the displaying subsystem (GUI). The global structure of the system. Early sketch based image retrieval systems were typically driven by queries comprising blobs of color or predefined texture. Later systems explored shape descriptors and spectral descriptors such as wavelets. Newly developed approach introduced a grid based approach to shape retrieval, categorizing the image into regular grids and locate photos using sketched doodle of object shape. Descriptors from each cell were concatenated to form a global image feature. However this offers limited invariance to changes in position, scale or orientation. A depiction invariant descriptor that encapsulates local spatial structure in the sketch and facilitates efficient codebook based retrieval was proposed by Hu et al.

#### A. Two Stage Architecture

To efficiently and effectively discover deep web data sources, Smart Crawling system is developed with two stage architecture, site locating and in-site exploring. The first stage finds the most relevant site for a searchable topic, and then the second stage finds searchable forms from the site. Specifically, the site locating stage begins with a seed set of sites in a site database. Seeds sites are candidate sites given for Smart Crawling system to start crawling, which starts with following URLs from chosen seed sites to explore other pages and domains. When there are number of unvisited URLs in the database is less than a threshold during the crawling, Smart Crawling system will performs "reverse searching" of known deep web sites for centre pages (highly ranked pages that have multiples links to other

domains) and feeds these pages back to the site database. Site Frontier extracts homepage URLs from the site database, which are ranked by Site Ranker to prioritize highly relevant sites. Adaptive Site Learner responsible for improvement of site ranker during crawling, which adaptively learns from features of deep-web sites (web sites containing one or more searchable forms) found. To achieve more accuracy among all results for a focused crawl, Site Classifier classified URLs into relevant and irrelevant for a user's given topic according to the homepage content. After finding the most relevant sites in the first stage, the second stage will performs efficient in-site exploration for excavating searchable forms. Link Frontier stores links of sites and corresponding pages are fetched and embedded forms are classified by Form Classifier to find searchable forms.

#### B. Site Locating

This stage finds relevant sites for a user's topic, consisting of site collecting, site ranking, and site classification.

#### C. Site Collecting

The previous crawler follows all newly found links. In contrast, our *Smart Crawling system* strives to minimize the number of visited URLs, and at the same time it increases *the* number of deep websites. To achieve these goals, using the links in downloaded WebPages is not sufficient. This is because a website usually contains a small number of links to other sites, even for some large sites. For instance, amazon.com contains 54 such links out of a total of 500 links thus, finding out-of-site links from visited WebPages may not be enough for the Site Frontier.

### CONCLUSION

We proposed the spatial constraint and coherent constraint to filter the false matches; this will validate the effectiveness of our framework. Though our method achieves significant performance improvement, demand of multimedia applications are increasing over the internet. Thus, the importance of image retrieval and image mining has increased. The proposed scheme can be considered competitive candidate in color image retrieval application. We pay attention to the area where the two clusters come closer to each other by applying hierarchical clustering algorithm to featured database.

### ACKNOWLEDGMENT

We have immense pleasure in presenting the paper "Smart Crawling System, A Review". We would like to thank MET's Institute of Engineering, Nashik for providing all the required facilities.

### REFERENCE

1. Tianjia Shao, Weiwei Xu, Kangkang Yin, Jingdong Wang, Kun Zhou, Baining Guo, "Discriminative Sketch-based 3D Model Retrieval via Robust Shape Matching", Volume 30, 2011.
2. Huda A. Abdulbaqi, Ghazali Sulong, Soukaena Hassan Hashem, "A Sketch Based Image Retrieval: A Review of Literature", Journal of Theoretical and Applied Information Technology, Vol. 63 No.1, 10 May 2014.
3. Radhika Bairagade, Nirmala Singh, Nikita Afre, Durga Bhamare, "A Survey on Smart Crawler: A Deep-Web Harvesting Approach", International Journal of Advanced Research in Computer Engineering & Technology (IJARCET), Volume 4 Issue 9, September 2015.
4. Neetesh Prajapati, G.S.Prajapti (Research Guide & HOD), and "Sketch Based Image Retrieval System for the Web - A Survey", IJCSIT: Vol. 6, 2015.
5. Sneha A. Ghumatkar, Prof. Archana C. Lomte, Prof. Gayatri Bhandari (Computer Department, JSPM'S ), "A Survey Paper on Web Crawler", International Journal of Advance Foundation and Research in Computer (IJAFRC) Volume: 2, Issue 12, ISSN 2348 – 4853, December - 2015.
6. Ms. Asmita A. Desai, Prof. Mrs. Aparna S. Shinde, Dr. Mrs. P. Malathi, "Sketch Based Image Retrieval System", IJARCET Volume 3 Issue 9, September 2014.
7. Prachi A. Gaidhani, S. B. Bagal, "Survey paper on Sketch Based and Content Based Image Retrieval", HAL Id: hal-01256024, 14 Jan 2016.
8. Nimisha Jain, Pragya Sharma, Saloni Poddar, Shikha Rani, "Smart Web Crawler to Harvest the Invisible Web World", IJIRCC: Vol. 4, Issue 4, and April 2016.
9. Archana Pondkule, Shital Khomane, Vaibhav Taware, "A Two-stage Smart Crawler : Harvesting Deep-Web Interfaces", DOI 10.4010/2016.571, Volume: 6.
10. Ms. Asmita D Rathod, "SMART CRAWLER: A TWO-STAGE CRAWLER FOR EFFICIENTLY HARVESTING DEEP-WEB INTERFACES", International Journal Of Innovation In Engineering Research And Technology [IJERT], ISSN: 2394-3696, Volume 3, APR.-2016.
11. Shu Wang, Jian Zhang, "Sketched-Based Image Retrieval Through Hypothesis-Driven Object Boundary Selection with HLR Descriptor", IEEE-vol.17, No. 7 July 2015.
12. Feng Zhao, Chang Nie, Hai Jin, "Smart Crawler: A Two-Stage Crawling For Efficiently Harvesting Deep-Web Interfaces", IEEE trans on Services Computing Volume: 99, PP Year: 2015.