

**INTERNATIONAL JOURNAL OF ENGINEERING SCIENCES & RESEARCH
TECHNOLOGY**

**CONTENT BASED INFORMATION RETRIEVAL FOR DIGITAL LIBRARY USING
DOCUMENT IMAGE**

Roshni S. Tadse*, L. H. Patil, C. U. Chauhan

* Research Scholar, Department of Computer Science and Engineering, Priyadarshini Institute of Engineering and Technology Nagpur (MS), India
Asst.Professor, Department of Computer Science and Engineering, Priyadarshini Institute of Engineering and Technology Nagpur (MS), India
Asst.Professor, Department of Computer Science and Engineering, Priyadarshini Institute of Engineering and Technology Nagpur (MS), India

DOI:

ABSTRACT

In the recent year, the using of mobile devices has perceive an emerging need for improving the user experience of digital library for search, with various applications such as education, location search and product retrieval, There simply compare the query to the databases images; those are match that images are retrieve from the database, searching and response time of delivery staying a challenging issues in mobile document search previously lots of work has been done on search engine, retrieving the document from the database without analyzed the image. In The proposed method, Information retrieval for image based query automatically with a mobile document information retrieval framework, consisting of a FP-growth is proposed finding frequent pattern from the retrieve document to optimize the result.

KEYWORDS: Digital library, FP-growth, Information Extraction, Mobile device, keyword Extraction.

INTRODUCTION

As this limits the applicability of search engines (images that do not coincide with textual data cannot be retrieved), thus developing a methods that generate description words for a picture automatically. Although keyword-based indexing techniques are popular and the method of choice for image retrieval engines. A method that generates such descriptions automatically could therefore improve image retrieval by supporting longer and more targeted queries, by creating as a short description of words for image's content, and by using the question-answer interfaces. The mobile devices has witnessed an emerging need to improve the user experience of digital library browsing and search, with various applications such as education, augmented reality, location search and product retrieval.

Data mining, the extraction of hidden predictive information from large databases, is a powerful new technology with great potential to help companies focus on the most important information in their data warehouses. Data mining tools predict future trends and behaviors, allowing businesses to make proactive, knowledge-driven decisions.

The automated, prospective analyses offered by data mining move beyond the analyses of past events provided by retrospective tools typical of decision support systems. Data mining tools can answer business questions that traditionally were too time consuming to resolve. They scour databases for hidden patterns, finding predictive information that experts may miss because it lies outside their expectations.

Most companies already collect and refine massive quantities of data. Data mining techniques can be implemented rapidly on existing software and hardware platforms to enhance the value of existing information resources, and can be integrated with new products and systems as they are brought on-line. When implemented on high performance

client/server or parallel processing computers, data mining tools can analyze massive databases to deliver answers to questions such as, "Which clients are most likely to respond to my next promotional mailing, and why?" In the recent years, digital library has played an important role in accessing the corpus of massive scanned documents stored in the digital image format. Content based retrieval could be a promising solution to facilitate pervasive and efficient access of the document images. In the typical scenario, a query is formulated as a photo that captures the visual objects of user interest, for example, a book cover, a document page. The visual query is sent to the server end, where the visually similar documents are matched and returned. To improve the image matching efficiency, the extracted visual signatures of database images have to be indexed, typically by an inverted indexing table. Compared to typing query keywords, a snapped photo based query undoubtedly simplifies the input of a user query. Furthermore, in some specialized domains like searching ancient hieroglyph, content based queries retain as the effective approach. The recent proliferation of mobile devices has witnessed an emerging need to improve the user experience of digital library browsing and search, with various applications such as education, augmented reality, location search and product retrieval.

RELATED WORK

[1] With the rapid multiplication use of mobile devices, previous years have find mobile searching techniques into digital library. Such a electronic device application's output has introduced the unique challenges in text document image search. The mobile photograph prefer hard to extract features from particular object's region of documents. In addition, searching and response time of query delivery bring out the challenging part in mobile document search. In this paper, In this paper propose a framework that is novel mobile document image retrieval framework, which is having the robust Local Inner-distance Shape Context (LISC) descriptor of line drawings, a Hamming distance KD-Tree for scalable, as well as a JBIG2 based query compression scheme with an OTSU based binarization, to reduce the response time of query delivery which maintaining query quality in terms of search performance.

Content based image retrieval (CBIR) [3] is the task of searching digital images from a huge database basis on the extraction of features, like as color, texture, shape of the image. Almost the research has been in CBIR which has been carried out with whole queries which were present in the database. This paper conclude the usefulness of CBIR techniques for retrieval of incomplete and misrepresented queries. [4] Text mining with information extraction having a two techniques BWI and RAPIER boosted wrapper induction and Robust automated production of IE rules. In this paper having a framework for text mining where combining the information extraction and knowledge data discover for text mining. It proposed for IE enabling the application of KDD to unstructured text corpora. KDD can discover predictive rules for improving IE performance. It is critical to the development of effective text mining systems for computational linguistics and machine learning. Open Language Learning for information extraction which extract relational tuples from text. Where OLLIE algorithm is implement to addresses the limitation of open information extraction and increase the precision.

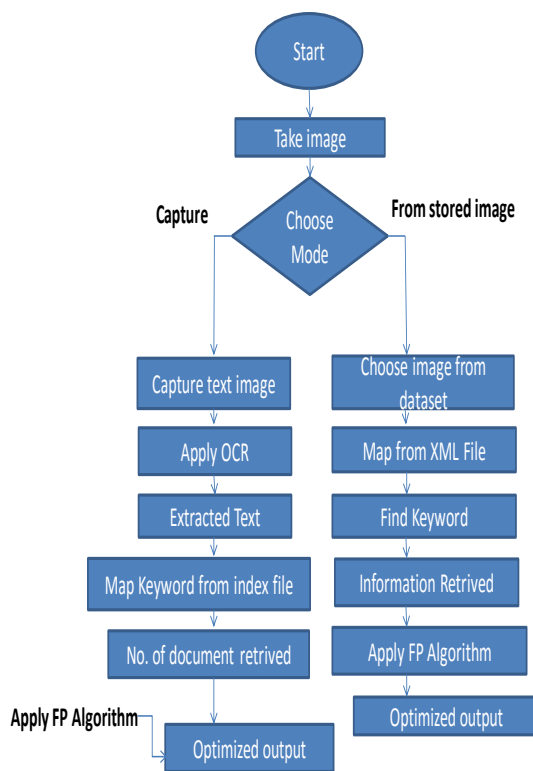


Figure 1: FlowChart of Proposed Work

In the above flowchart shown the complete scenario of this paper where the figure1,attemp the capture mode and stored image from the dataset, flow of the paper’s scenario. These two mode has been decided by the user after this stage there are two way for extracting the information first is by applying OCR and the second one is mapping through the XML . After selecting the choice mode of selecting or capturing the image the further work will done on server side.

Information Retrieval is based on query you specified what information you need and it is return in human understandable form and information extraction is about structuring unstructured information given some sources of all the information is structured in a form that will be easy for processing this will not necessary be in human understandable form. It can be only use for computer programs. In the left part of the flow chart shown that after applying the OCR the text is extracted. First the program analyze the structure of document image. It divides page into elements such as blocks of texts,tables images etc. The are divided into words and words into characters after processing the huge number of such probabilistic hypotheses, the program finally takes the decision presenting you the recognized text. With dictionary support the program ensures even more accurate analysis and recognition of document. Map the keyword from indexed file where already indexing is done into the database now, number of files are retrieved after recognizing the text. If the input query is select from dataset images then it must be mapped already in backend. The text here can map the attributes of java object to a combination of xml simple and complex types using a wide variety XML mapping type. From both side the number documents are retrieved. For getting the optimized output the FP-growth has been applied for retrieving a meaningful document.

Graphical User Interface:

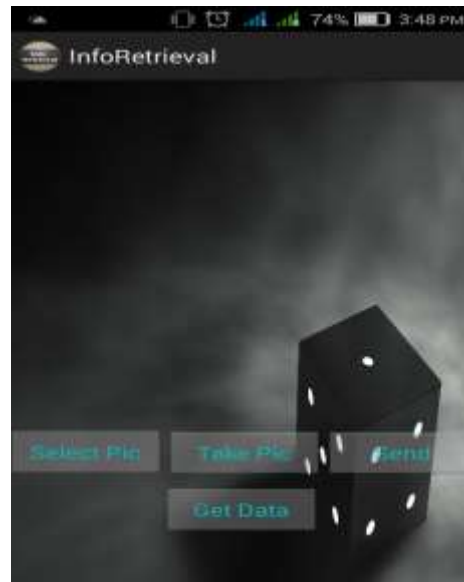


Figure2 Graphical User Interface

In the above image shown the graphical user interface which Information Retrieval Apk(application) for according to the so far idea. This Apk is comes under the user side that is in user's mobile. Where user can easily access the front view for sending the query to the server. The image is formulated as query and send it to server, where the two options for the user one capture the image and sending to the server and second one is select the image from the dataset for formulated the query to send the server.

In the above image having three buttons that is select picture, take picture and send which is easily handled by the user or non technical person too. Generally, GUI must in that way which is easily handled by navy user .



Figure 3: Capturing the image by mobile device

In the proposed research plan of work shows that the visual query is captured by the user and sent to the server end, where the visually similar documents are matched and returned. Many of the search engines deployed on the web retrieve images without analysing content present in the image; they match user queries against collocated textual information.

Firstly, it take a picture or capture the image through mobile which should be a text image. After that it goes into the choose mode onto the mobile. Here there has an option to take as an input whether it is captured image or from stored image of the dataset. Hence, the choose mode will help to select the mode of the input image easily.



Figure 4: Retrived File based on send query

The main objective of this work is for indexing text data and using this technology to implement searching and indexing of 100 million XML files from the CoPhIR (Content-based Photo Image Retrieval) data set. TheCoPhIR data set is a test collection that serves as the basis of the experiments on content-based image retrieval techniques. The data collected within represents the world largest multimedia metadata collection available for research purposes, containing visual and textual information regarding over 100 million images. Organizes 100 million image collections based on the Content based information retrieval. Indexing and searching mechanisms are based on the concept of structured peer-to-peer networks, which makes its approach highly scalable and independent of the specific hardware

infrastructure on which it runs. The resulting index has to be optimized for searching in descriptive image data. Before the actual indexing, data cleansing has to be executed. Subsequently, an index can be built and optimized for user searches



Figure 4: Optimized output

CONCLUSION

The previous approach was dealing with simple heuristic rules which did not give proper results. We will propose a system to find frequent patterns using FP-Growth algorithm which gave efficient results compared to previous system. These frequent set of keywords improved the accuracy of document retrieval from digital library using mobile images. We have retrieved all the related information from our digital library which will be useful for finding frequent patterns and ultimately retrieving meaningful documents.

The proposed method infer that mobile document images are sending to the server and retrieved the document of that particular image. An efficient method of retrieving document images from content based information retrieval for document and searching methods that search the information with respect to content of images. Here getting the optimized output at the end user in an efficient manner. JPEG compression is there, to low complexity is introduced to reduce the query delivery latency while maintaining comparable search accuracy.

REFERENCES

- [1] Ling-Yu Duan, Rongrong Ji "Towards Mobile Document Image Retrieval for Digital Library" IEEE Transactions on multimedia, vol. 16, no. 2, february 2014
- [2] Gulfishan Firdose Ahmed, Raju Barskar "A Study on Different Image Retrieval Techniques in Image Processing" International Journal of Soft Computing and Engineering (IJSCE) ISSN: 2231-2307, Volume-1, Issue-4, September 2011
- [3] Bikesh Kumar Singh, A. S. Thoke "Image information retrieval from incomplete queries using color and shape features", Signal & Image Processing : An International Journal (SIPIJ) Vol.2, No.4, December 2011
- [4] Raymond J. Mooney, Un Yong Nahm "Text Mining With Information Extraction" 4th International MIDP Colloquium, September 2003.
- [5] Mausam, Michael Schmitz, Robert Bart, Stephen Soderland, and Oren Etzioni "Turing Center University of Washington", Seattle.
- [6] R. Basri, L. Costa, D. Geiger, and D. Jacobs, "Determining the similarity of deformable shapes," Vision Res., vol. 38, no. 15, pp. 2365–2385, 1998.
- [7] A. Ratarangsi and R. T. Chin, "Scale-based detection of corners of planar curves," in Proc. IEEE 10th Int. Conf. Pattern Recognition, 1990, vol. 1, pp. 923–930.
- [8] Z. Tu and A. L. Yuille, "Shape matching and recognition-using generative models and informative features," in Proc. ECCV, 2004, pp. 195–209.

- [9] X. Chen, H. Nelson, and H. Yung, "Corner detector based on global and local curvature properties," *Opt. Eng.*, vol. 47, no. 5, p. 057008, 2008.
- [10] J. Pu and K. Ramani, "On visual similarity based 2D drawing retrieval," *Comput. Aided Design*, vol. 38, no. 3, pp. 249–259, 2006.
- [11] D. Chen, S. Tsai, V. Chandrasekhar, G. Takacs, J. Singh, and B. Girod, "Tree histogram coding for mobile image matching," in *Proc. DCC*, 2009, pp. 143–152.
- [12] V. Chandrasekhar, G. Takacs, D. Chen, S. Tsai, R. Grzeszczuk, and B. Girod, "CHoG: Compressed histogram of gradients a low bit-rate feature descriptor," in *Proc. CVPR*, 2009, pp. 2504–2511.
- [13] D. Lowe, "Distinctive image features form scale-invariant keypoints," *Int. J. Comput. Vision*, vol. 60, no. 2, pp. 91–110, 2004.
- [14] H. Bay, T. Tuytelaars, and L. Van Gool, "SURF: Speeded up robust features," in *Proc. ECCV*, 2006, pp. 404–417.
- [15] T. K. Bhowmik, U. Bhattacharya, and S. Parui, "Recognition of bangla handwritten characters using an MLP classifier based on stroke features," in *Neural Information Processing*. Berlin, Germany: Springer, 2004, pp. 814–819.
- [16] R. Ji, L.-Y. Duan, J. Chen, H. Yao, J. Yuan, Y. Rui, and W. Gao, "Location discriminative vocabulary coding for mobile landmark search," *Int. J. Comput. Vision*, vol. 96, no. 3, pp. 290–314, 2012.
- [17] W. Zhang and J. Kosecka, "Image based localization in urban environments," in *Proc. 3DPVT*, 2006, pp. 33–40.
- [18] G. Schindler and M. Brown, "City-scale location recognition," in *Proc. CVPR*, 2007, pp. 1–7.
- [19] E. Eade and T. Drummond, "Unified loop closing and recovery for real time monocular SLAM," in *Proc. BMVC*, 2008, vol. 13, p. 136.
- [20] T. Yeh, K. Tollmar, and T. Darrell, "Searching the web with mobile images for location recognition," in *Proc. CVPR*, 2004, vol. 2, pp. II–76.
- [21] V. Santosh, D. D'Souza, T. Kavitha, and J. Radhakrishnan, "Randomly-oriented k-d trees adapt to intrinsic dimension," in *Proc. ARCS Annu. Conf. Foundations of Software Technology and Theoretical Computer Science*, 2012, pp. 48–57.
- [22] Paolo Bolettieri, Andrea Esuli, Fabrizio Falchi, Claudio Lucchese, Raffaele Perego, Tommas Piccioli, Fausto Rabitti CoPhIR: a test collection for content-based image retrieval. *CoRR*, vol. abs/0905.4627, 2009 ,