

# Facilitating Document Annotation using content and Querying value-A Review

Miss. Sravya<sup>1</sup>, Mr. Laxman<sup>2</sup>, Mr. Sravan<sup>3</sup>, Miss. Sravanthi<sup>4</sup>  
B.Tech Student Department of CSE,

Balaji Institute of Technology & Sciences, Narsampet, Warangal, Telangana, India.

## ABSTRACT:

Annotation is an effective information retrieval which is used to add additional information to the documents. Annotation helps to analyze and retrieve the required documents easily based on the query given by the user. Annotation can be applied to several types like images, documents etc. For annotation, firstly related attributes have to be specified. For making explication, people need to work very hard to read the content of the document thoroughly to find which statements has to be annotated. For that in this paper two performances of annotations are discussed and detail on the proposed system is explained. Firstly, CADs system is used to allow the effortless sharing of documents and at the same time serving semi structured queries in the database. LableMe is also used for the comment of image annotation. These techniques help the data or images to retrieve quickly and accurately than the existing system. The above Methods used in this paper are best for the retrieving the related documents.

**Keywords:** Document Annotation, CADs System, Image Annotation, LableMe Annotation Tool.

## 1. INTRODUCTION:

Data mining is the process of searching large amount of data which is used to identify the related patterns .The main aim of data mining is used to retrieve the large amount of data where that data is converted into required patterns .Information extraction is the method used to extract information from large number of documents .Annotation and content extraction are used for extraction of information other multimedia processing methods can also be used Annotations provides users more suitable data by violating the unrelated information and to understand the document effectively and this process increase the efficiency of the searching where accurate results will be provided based on the user query.

Text annotations has a variety of applications some of the applications include social reading, writing and educational applications. Annotations are the authors opinion to find the data ,it is the more powerful method .Making annotations is hard because people has to read the entire document carefully to identify which part has to be annotated. Automatic annotation is used for saving of time and makes the documents in structured format .In addition to this, image annotation assigns a Meta data automatically Documents use annotations for saving main attributes in the database for further searching. Advantages of automatic annotation are:

-fast access

-Automated annotation methods generally gives more annotations compared to manual.

For annotation, two methods has to be used there are CADs and LableMe. The main purpose of the CADs system is used to lower the cost of annotated document that should be used for the given queries. CADs application is used in Business Continuity Information Network (BCIN) of southflorida. CADs system searches the content of the document and create a insertion form which is adaptive which contains the needed information and required best attributes of the document along with the attribute values of the given document text.

## 2 DOCUMENT ANNOTATION PROCESS:

we will design the methodology for the documents as well as the images also that will reduce the time and space along with multiple search .The proposed system , will provide the image annotation along with the document that are used to find effective information in less time. The two annotations techniques implemented are:

### 2.1 CADs System:

The main aim of CADs is reduce the cost that is commonly used for semi structured queries. Fig (1) represents the CAD workflow it consists of two types of actors they are: producers and consumers. Producers are used to upload the data by using interactive insertion form in CADs, where as consumers search for the relevant document using adaptive query forms. Its objective is to create structured annotated document with less cost.

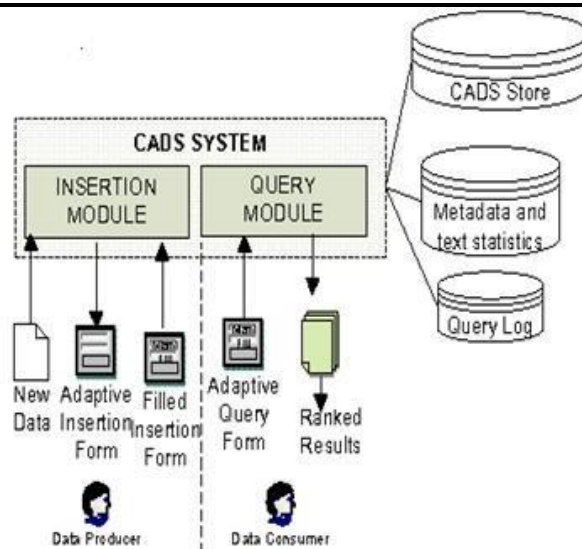


Figure: 1 CADS workflow

The CADS system consists of two modules: Insertion module And Query module. In the insertion module the document Submission is done and then CADS analyzes the document and Creates an adaptive insertion form it contains the attribute values to annotate the document. The user fills the require data and submit it to the database. In the same way, in the query phase the user interacts with the adaptive query form where it contains some attributes if the user need to add some more attributes that is also available and it also contains description attribute where the description of the document can be added. Finally, the CADS system will give the most important data Pieces and also it returns the whole document.



### 3 LITERATURE SURVEY:

[1] Eduardo J. Ruiz, Vagelis Hristidis, and Panagiotis G. Ipeirotis : Facilitating Document Annotation Using Content and Querying Value. That paper [1] presents the algorithms that identify structured attributes that are likely to appear within the document, by jointly utilizing the content of the text and the query workload. [2] CADS: That paper [2] proposed CADS system, which are used as a Collaborative Adaptive Data Sharing platform, and are a data sharing platform where the integration and annotation take place at the time of data insertion i.e. production and querying i.e. consumption actions. A main aim of CADS is to influence the information demand for creation of adaptive insertion and query forms. [3] Proximity Ranking: The Recent studies show that the term proximity are highly correlated with relevancy of document, and proximity aware ranking increases the top results precision significantly. And, there are only few studies which increase proximity-aware searching query efficiency using techniques of early termination [3], [4]. The techniques which are discussed in [3], [4] generate an additional inverted index for each term pair, which results in a large space. [4] studied only the problem for queries with two keywords. [4] LableMe: B. Russell, A. Torralba, K. Murphy, and W. Freeman : propose a paper "Label Me: A Database and Web Based Tool for Image Annotation". A tag prediction for images is proposed in that paper [5]. It proposes web-based tool for easy image annotation and instant sharing of annotations. It detects the objects and finds similarity with existing dataset. It helps for image search in web. [5] A tag prediction for images is proposed in that paper [5]. It proposes web-based tool for easy image annotation and instant sharing of annotations. [6] Instant Search: The integration of proximity information in instant fuzzy search for achieving the better complexities is explained in many recent studies focused on the instant search. The studies in [6] proposed query and indexing techniques to support the instant search. Li et al. [6] studied the instant search on relational data which is modelled as a graph.

### 4 SYSTEM OVERVIEW:

The description of the system is as follows as shown in the Architecture of the document annotation, the user uploads the documents in the document uploader that is already Annotated where the CAD system retrieve and calculate them in the database using CADS. The user can upload the images in the image uploader using LableMe annotation tool that are stored in the database suppose if the user searches for a file that may a filename, label, content, query or image name too. An algorithm is used to check the queries by managing the history of users based on that it will return the related image or document.

### 5 LableMe Annotation Tool:

LableMe is the annotation tool and a database. It provides Drawing functionalities and sharing of the images and also provides the quality of labelling. The main goal of this tool is to provide drawing interface in any kind of platforms. The User can label and select the image. Labelling of image can be done by clicking the control boundaries of the object after finishing the dialog button will be displayed for object name. This is saved in the Label Me database and corresponding image Will be displayed which is available for download and for the view of the user. This tool is very easy and simple to use.

### 6 Information Extraction Algorithm:

This algorithm is used extract the contents of a text file, the below figure shows how extraction takes place Content value and query value for the attributes, If the attribute value is highest it will generate the insertion from the documents annotated are stored

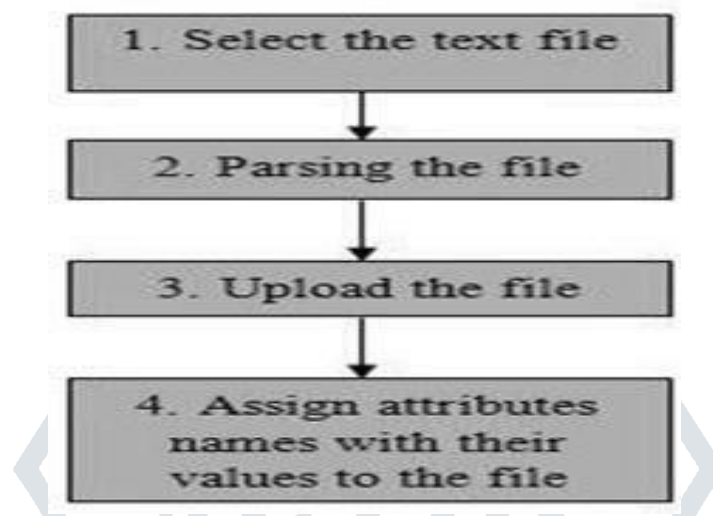


Figure (2) algorithm flow

Step 1: Select a text file for extraction.

Step 2: parse a text file. Ignore stop words from it and count frequency of high querying key words appearing in a single document.

Step 3: upload the file on server.

Step 4: then fill all the annotations which are relevant to the document which can be used for query based searching.

**Example:** year=2012, location="Nasik", author="Bill Gates" etc.

### 7 QV, CV Computation and Combining: Algorithm:

1. Enter the queries for retrieving the document
2. Split the queries and pass to the database for retrieving.
3. Check the results and display the results for the user.
4. For more efficient and effective results, user has to enter the maximum queries.

The below figure shows that a publisher has to choose the file for uploading, the uploaded document is parsed where the stop words are ignored and key words are stored in the database

Enter Category Name

---

Select Category

Select TextFile  imp.txt

Figure (3)



Figure (4)



Figure 5:

After uploading the file by a publisher an attribute name, domain and respective values will be assigned. The reviewer enters the word “facilitating” and all the documents containing the word will be shown.

Query Name	Query Value	Query Type
year	2012	Number
author	suresh	Text
year	2001	Number
author	Akshay	Text
year	2000	Number
loc	nashik,pune,mumbai	Text
loc	mumbai	Text
Title	Supporting Efficient and Scalable Multicasting over Mobile Ad Hoc Networks	Text
loc	pune	Text
loc	mumbai	Text

Figure (6)

The list of attributes/annotations with their respective values and types can be seen by the end user which is used for query based searching.

## 8 RESULT AND DISCUSSION:

This system helps to generate the results automatically for the documents as well as the images that are uploaded by user by using the CADS system and LableMe annotation tool. the user will retrieve the ranked documents of his mentioned queries within the less time .

## 9 CONCLUSION:

Now-a-days, the annotations can generate suggestions based on the given paper but in future it can generate the more information such as relevant documents and additional information about the annotation. Annotation technique is more powerful tool used for retrieval of the documents or images in future, document clustering can be used This paper explains the new approaches of the related documents retrieval by using multiple annotation method or by searching and ranking techniques. This system used to increase the document visibility and efficient querying of the user. Query expansion algorithm is used to remove stop and stem words. This system is used to provide efficient and fast retrieval of document with less time and space

## 10 REFERENCES:

- [1] Ahamed, B. B., & Hariharan, S. (2012). Implementation of Network Level Security Process through Stepping Stones by Watermarking Methodology. *International Journal of Future Generation Communication and Networking*, 5(4), 123-130.
- [2] H. Yan, J. Wen, S. Shi, F. Zhang, T. Suel, "Efficient term proximity search with the term-pair indexes," *CIKM*, 2010, pp. 1229-1238.
- [3] H. Bast, A. Chitea, F. Suchanek, Weber, "Ester : efficient search on text, entities, and relations," *SIGIR*, 2007.
- [4] B. Russell, A. Torralba, K. Murphy, and W. Freeman: propose a paper "LabelMe: A Database and Web-Based Tool for Image Annotation".
- [5] Md. Abu Nisar Masud, Md. Munasir Mamun, "A General Approach to Natural Language Generation" In *Proceeding of IEEE, INMIC, 2003*.
- [6] Vishal A Patil, & Pankaj Khambre. (2015). Survey on Facilitating Document Annotation using Content and Querying Value. *International Journal of Science and Research*. Vol. 4.
- [7] Eduardo J. Ruiz, Vagelis Hristidis, & Panagiotis G. Iperiotis. (2014). Facilitating Document Annotation using Content and Querying Value. *IEEE Transaction on Knowledge and Data Engineering*.
- [8] Microsoft Sharepoint. Available at <http://www.microsoft.com/sharepoint/>
- [9] Google Base. Available at <http://www.google.com/base>
- [10] MJ Cafarella, J Madhavan, & A Halevy. (2009). Web-scale extraction of structured data. *SIGMOD Rec*. Vol. 37, pp. 55-61.
- [11] SR Jeffery, MJ Franklin, & AY Halevy. (2008). Pay-as-You-Go User Feedback for Data space Systems. *Proc. ACM SIGMOD Int'l Conf. Management Data*.
- [12] CD Manning, P Raghavan, & H Schütze. (2008). *Introduction to Information Retrieval*. Cambridge University Press.
- [13] J Madhavan et al. (2007). Web-Scale Data Integration: You Can Only Afford to Pay as You Go. *Proc. 3<sup>rd</sup> Biennial Conf. on Innovative Data Systems Research (CIDR)*.
- [14] Gomathi, M., & Ahamed, B. B. Socio-Technical Accordance Perspective For Software Implementation Correlation With Fault Aptitude.
- [15] JM Ponte, & WB Croft. (1998). A Language Modelling Approach to the Information Retrieval. *Proc. 21st Ann. Int'l ACM SIGIR Conf. Research and Development in Information Retrieval (SIGIR '98)*, pp. 275-281. <http://doi.acm.org/10.1145/290941.291008>
- [16] A Halevy, Z Ives, D Suciu, & I Tatarinov. (2003). Schema Mediation in Peer Data Management Systems. *Proc. 19<sup>th</sup> Int'l Conf. Data Eng.*, pp. 505-516.
- [17] Punidha, R., avithra K, Swathika R, and Sivaram M, "Preserving DDoS Attacks sing Node Blocking Algorithm." *International Journal of Pure and Applied Mathematics*, Vol.119, o. 15, 2018, pp 633-640. <https://acadpubl.eu/hub/2018-119-15/3/473.pdf>
- [18] Ahamed, B. B., & Ramkumar, T. (2015). Deduce User Search Progression with Feedback Session. *Advances in Systems Science and Applications*, 15(4), 366-383.

**11. AUTHORS BIBLIOGRAPHY:****SRAVYA MERUGU**

Pursuing B.Tech in Computer Science and Engineering, Balaji Institute of Technology & Science, Warangal, Telangana, India.

**LAXMAN KALLEPU**

Pursuing B.Tech in Computer Science and Engineering, Balaji Institute of Technology & Science, Warangal, Telangana, India.

**SRAVAN KADUDHULA**

Pursuing B.Tech in Computer Science and Engineering, Balaji Institute of Technology & Science, Warangal, Telangana, India

**SRAVANTHI DHUSHANPALLY**

Pursuing B.Tech in Computer Science and Engineering, Balaji Institute of Technology & Science, Warangal, Telangana, India.

