



# INTERNATIONAL JOURNAL OF PURE AND APPLIED RESEARCH IN ENGINEERING AND TECHNOLOGY

A PATH FOR HORIZING YOUR INNOVATIVE WORK

## REVIEW PAPER ON IMPLEMENTATION OF DOCUMENT ANNOTATION USING CONTENT AND QUERYING VALUE

MISS. ANUPAMA V. ZAKARDE<sup>1</sup>, DR. H. R. DESHMUKH<sup>2</sup>

1. Dept. of Computer Science, I.B.S.S. College of Engineering, Amravati

2. Prof. & Head IBSS college of Engineering, Amravati.

Accepted Date: 05/03/2015; Published Date: 01/05/2015

**Abstract:** Annotation plays a major role in a user's reading of a document: from elementary school students making notes on text books to professors marking up their latest research papers. A common place for annotations to appear is in the margin of a document. Surprisingly, there is little systematic knowledge of how, why and when annotations are written in margins or over the main text. This project investigates how margin size impacts the ease with which documents can be annotated, and user annotation behavior. The research comprises of a two part investigation: first, a survey which examines margins and their use in physical documents; secondly, we evaluate document reader software that supports an extended margin for annotation in digital documents. This work present a novel alternative approach that facilitates the generation of the structured metadata by identifying documents that are likely to contain information of interest and this information is going to be subsequently useful for querying the database. The approach relies on the idea that humans are more likely to add the necessary metadata during creation time, if prompted by the interface; or that it is much easier for humans (and/or algorithms) to identify the metadata when such information actually exists in the document, instead of naively prompting users to fill in forms with information that is not available in the document.

**Keywords:** Document annotation, adaptive forms, collaborative platforms

Corresponding Author: MISS. ANUPAMA V. ZAKARDE



PAPER-QR CODE

Access Online On:

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

How to Cite This Article:

Anupama V. Zakarde, IJPRET, 2015; Volume 3 (9): 821-826

## INTRODUCTION

During several degree programs of universities academic institutes we have research component of varying duration from six months to four-five years. As a first activity in the research, students have to survey literature related to their domain of interest to define their proposed activity. They collect research papers and other publication either from web sites of professional societies like IEEE, ACM, and LNCS or from printed copy of journals available in their library. While going through these research publications, they marks underline on some important parts, or they highlights some words or phrases or whole sentences or paragraph. They also write their notes, observations, remarks, questions etc either on the same document or on the separate sheet of paper. These highlights or underline or comments/observation may be about entire paper or part of them. At some point of time, they collate and integrate these observations to identify and define their research problems. On completions of their degree programs these knowledge represented in the form of observation and thoughts are lost as they are not saved and shared by the next batch of students. These observation/comments/highlights/underline are very valuable knowledge resource not only for the current reader but also for future generation of students who are likely to work in the same area. However, at present these knowledge resources are not available to future generation as they are not available in electronic form and are not sharable. This work is motivated by desire to provide a tool which provides a facility to record their comments, notes, observation, and explanation, highlights, underline etc. either on document or on another comments and evaluate the collective sentiments of the researchers over the document. These collective sentiments of annotators may be used as an indicator of quality or usefulness of the documents. In this project, research that demonstrates the significance of margin space in annotating both digital and physical documents. Though the issue of margin space may seem trivial, there is a lack of concrete research across much of the field of annotation. Provide detailed evidence on chosen topic and demonstrate how digital document reader software can be significantly improved by changing their interaction design, informed by observation of actual user behavior.

**LITERATURE REVIEW & RELATED WORK:** Jain, P.G. Ipeirotis [1] Introduced a rigorous model for estimating the quality of the output of an information extraction system when paired with a document retrieval strategy. How to generate a ROC curve that can generate a statistically robust performance characterization of an extraction system, and then built statistical models that use the ROC curves concept to build the *quality curves* that predict the performance of

coupling an extraction system with a retrieval strategy. Analysis helps predict the execution time and output quality of an execution plan.

R.T. Clemen & R.T. Winkler [2] In this paper, variety of approaches for combining probabilities. The inquiry has been how these models relate to principles of unanimity and compromise. Those models that provide for the most general patterns of dependence among sources are the most complex in terms of their conformance to the principles.

B. C. Russell , A. Torralba, K. Murphy, W. Freeman [3] In this paper, A web-based image annotation tool that was used to label the identity of objects and where they occur in images. We collected a large number of high quality annotations, spanning many different object categories, for a large set of images, many of which are high resolution. Presented quantitative results of the dataset contents showing the quality, breadth, and depth of the dataset. We showed how to enhance and improve the quality of the dataset through the application of WordNet, heuristics to recover object parts and depth ordering, and training of an object detector using the collected labels to increase the dataset size from images returned by online search engines.

M.J. Cafarella, J. Madhavan, and A. Halevy [4] In this paper, there are three systems that perform information extraction in a domain-independent fashion, and therefore can applied to the entire Web. In all three cases, a side result of the extraction is a set of entities, relationships and schemata that can be used as building blocks for the Web knowledge base and for additional semantic services.

O. Etzioni, M. Banko, S. Soderland, and D.S. Weld [5] This paper introduces Open IE from the Web, an unsupervised extraction paradigm that eschews relation-specific extraction in favor of a single extraction pass over the corpus during which relations of interest are automatically discovered and efficiently stored. The paper also introduces TEXTRUNNER, a fully implemented Open IE system, and demonstrates its ability to extract massive amounts of high-quality information from a nine million Web page corpus.

V. Uren, P. Cimiano, J.Iria , S. Handschuh, M. V. Vera , E. Motta , F. Ciravegna [6] In this paper, documents are central to KM, but intelligent documents, created by semantic annotation, would bring the advantages of semantic search and interoperability. These systems need automation to support annotation, automation to support ontology maintenance, and automation to help maintain the consistency of documents, ontologies and annotations.

**ANALYSIS OF PROBLEM:** There are many application domains where users create and share information; for instance, news blogs, scientific networks, social networking groups, or disaster management networks. Current information sharing tools, like content management software (e.g., Microsoft Share-Point), allow users to share documents and annotate (tag) them in an ad hoc way. Similarly, Google Base allows users to define attributes for their objects or choose from predefined templates. This annotation process can facilitate subsequent information discovery. So there is requirement of an adaptive technique for automatically generating data input forms, for annotating unstructured textual documents, such that the utilization of the inserted data is maximized, given the user information needs and an algorithms to seamlessly integrate information from the query workload into the data annotation process, to generate metadata that are not just relevant to the annotated document, but also useful to the users querying the database.

**PROPOSED WORK & OBJECTIVES:** In this propose system Collaborative Adaptive Data Sharing platform, which is an “annotate-as-you-create” infrastructure that facilitates fielded data annotation. A key contribution of our system is the direct use of the query workload to direct the annotation process, in addition to examining the content of the document. In other words we are trying to prioritize the annotation of documents toward generating attribute values for attributes that are often used by querying users. In Fig. 1(a) we show a report extracted from the National Hurricane Center repository, describing the status of a hurricane event in 2008. The report gives the current storm location, wind speed, warnings, category, advisory identifier number, and the date it was disclosed. Even though this is a text document, it contains implicitly many attribute names and values, for example, (Storm Category, 3). If we had these values properly annotated (e.g., as in Fig. 1b), we could improve the quality of searching through the database. For instance, Fig. 1c shows three sample queries for which the report of Fig. 1a is a good answer and the lack of the appropriate annotations makes it hard to retrieve it and rank it properly. In propose system, whole work will divide in particular model as in first model we will collecting related document on which annotation will perform file format such as PDF document and text will be input to our propose system. Then attribute suggestion will be done on the input document.

```
ZCZC MIATCPAT2 ALL
TTAA00 KNHC DDHHMM
BULLETIN
HURRICANE GUSTAV INTERMEDIATE ADVISORY
NUMBER 31A
NWS TPC/NATIONAL HURRICANE CENTER MIAMI FL
AL072008
600 AM CDT MON SEP 01 2008

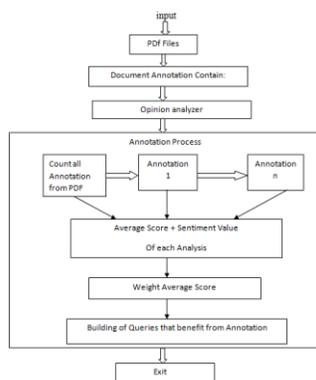
EYE OF GUSTAV NEARING THE LOUISIANA
COAST...HURRICANE FORCE WINDS OVER PORTIONS
OF SOUTHEASTERN LOUISIANA... A HURRICANE
WARNING REMAINS IN EFFECT FROM JUST EAST
OF HIGH ISLAND TEXAS EASTWARD TO THE
MISSISSIPPI-ALABAMA BORDER...INCLUDING THE
CITY OF NEW ORLEANS AND LAKE PONTCHARTRAIN.
PREPARATIONS TO PROTECT LIFE AND PROPERTY
SHOULD HAVE BEEN COMPLETED. A TROPICAL
STORM WARNING REMAINS IN EFFECT FROM
EAST OF THE MISSISSIPPI-ALABAMA BORDER TO
THE OCHLOCKONEE RIVER. GUSTAV IS MOVING
TOWARD THE NORTHWEST NEAR 16 MPH...26
KM/HR... ON THE FORECAST TRACK...THE CENTER
WILL CROSS THE LOUISIANA COAST BY MIDDAY
```

Figure(a):Example of unstructured Document

```
Storm Name = 'Gustav'
Storm Category = 3
Warnings = 'tropical storm'
```

Figure (b): Desirable annotations for the document above

Fig.1. Sample document and annotations.



In next module, proposed system will count number of annotation and suggest some desirable annotation for the above unstructured document. In next module, average score of each annotation will be calculated. So that desirable annotation will be properly shown as output. In finale module, according to the desirable annotation set of attribute queries will build that can benefit from the annotation. This project present two way to combine piece of explanation such as querying value and content value and it can suggest the attribute that improve the visibility document with respect to the query workload. The query workload might be improve annotation process and increase utility of shared data.

**Application:** A technical point of view, annotations are usually seen as metadata, as they give additional information about an existing piece of data. There are many application domains

where users create and share information; for instance, news blogs, scientific networks, social networking groups, or disaster management networks.

**Conclusion:** A large number of organizations today generate and share textual descriptions of their products, services, and actions. Such collections of textual data contain significant amount of structured information, which remains buried in the unstructured text. While information extraction algorithms facilitate the extraction of structured relations, they are often expensive and inaccurate, especially when operating on top of text that does not contain any instances of the targeted structured information.

## 7. References:

1. A. Jain and P.G. Ipeirotis, "A Quality-Aware Optimizer for Information Extraction," ACM Trans. Database Systems, vol. 34, article 5, 2009
2. R.T. Clemen and R.L. Winkler, "Unanimity and Compromise among Probability Forecasters," Management Science, vol.36, pp.767779, <http://portal.acm.org/citation.cfm?id=81610.81609>, July 1990.
3. B. Russell, A. Torralba, K. Murphy, and W. Freeman, "Label Me: A Database and Web-Based Tool for Image Annotation," Int'l J. Computer Vision, vol. 77, pp. 157-173, <http://dx.doi.org/10.1007/s11263-007-0090-8>, 2008, doi: 10.1007/s11263-007-0090-8.
4. M. J. Cafarella, J. Madhavan, and A. Halevy, "Web-Scale Extraction of Structured Data," SIGMOD Record, vol. 37, pp. 55-61, <http://doi.acm.org/10.1145/1519103.1519112>, Mar. 2009.
5. O. Etzioni, M. Banko, S. Soderland, and D.S. Weld, "Open Information Extraction from the Web," Comm. ACM, vol. 51, pp. 68-74, <http://doi.acm.org/10.1145/1409360.1409378>, Dec. 2008.
6. V. Uren, P. Cimiano, J.Iria, S. Handschuh, M. V. Vera, E. Motta, F. Ciravegna, "Semantic annotation for knowledge management: Requirements and a survey of the state of the art", 1570-8268/\$ doi:10.1016/j.websem.2005.10.002
7. Eduardo J. Ruiz, Vagelis Hristidis, and Panagiotis G. Ipeirotis, "Facilitating Document Annotation Using Content and Querying Value", IEEE transactions on knowledge and data engineering, vol. 26, no. 2, february 2014.