# ECADS: Annotation Query Technique using Dynamic Form Generation Model

Roopam Chaturvedi
M Tech (CSE)
LNCT Bhopal (M.P.) India

Vineet Richharya, PhD
HOD (CSE)
LNCT Bhopal (M.P.) India

Shweta Shrivastava
Assistant Professor (CSE)
LNCT Bhopal (M.P.) India

## ABSTRACT

A large number of systems today generate and share textual descriptions of their products, services, and actions. Such assemblage of textual data contain significant amount of structured information, which reside in the unstructured text. While related information extraction algorithms facilitate the extraction of structured relations, they are often costly and defective, especially when operating on top of text that does not contain any instances of the final targeted structured information. We present a good alternative approach that facilitates the generation of the structured metadata by identifying documents that are likely to have information of interest and this information is going to be subsequently useful for querying the database.

Our approach depends on the idea that instead of writing the query for each requirement it is easier to fetch them by using the annotated form based technique. In this approach user do not need to learn the correct query along with the previous information of the dataset. As a major contribution towards this paper, we present algorithms that identify structured attributes that are likely to find within the document, by jointly utilizing the content of the text and the query workload. Our experimental evaluation display that our approach generates superior results compared to approaches that go through textual content or only on the query workload, to identify attributes of our interest.

## Keywords

Data mining, ECADS, Annotation, Query Accessing, Content Value, Query Value

## 1. INTRODUCTION

In the recent years, a huge amount of data is being gathered and stored in databases everywhere across the globe, which is mainly coming from information industry and social networking sites. There is a need to extract and classify useful information and knowledge from such a data collected. Data mining is an interdisciplinary field of computer science and is referred to as extracting or mining knowledge from large databases. It is the process of performing automated extraction and generating the predictive information from a large database. It is actually the process of finding the hidden information or patterns from the repositories .The fields that use Data mining techniques include medical research, marketing, telecommunication, and stock markets, health care and so on.

This annotation process can facilitate subsequent information discovery. Many annotation systems allow only "untyped" keyword annotation: for instance, a user may annotate a weather report using a tag such as "Storm Category 3." Annotation strategies that use attribute-value pairs are generally more expressive, as they can contain more information than untyped approaches.

In such settings, the above information can be entered as (Storm Category, 3). A recent line of work toward using more expressive queries that leverage such annotations, is the "pay-as-you-go" querying strategy in Data spaces: In Data spaces, users provide data integration hints at query time. The assumption in such systems is that the data sources already contain structured information and the problem is to match the query attributes with the source attributes. Many systems, though, do not even have the basic "attribute-value" annotation that would make a "pay-as you- go" querying feasible. Annotations that use "attribute value" pairs require users to be more principled in their annotation efforts.

Users should know the underlying schema and field types to use; they should also know when to use each of these fields. With schemas that often have tens or even hundreds of available fields to fill, this task become complicated and cumbersome. This results in data entry users ignoring such annotation capabilities. Even if the system allows users to arbitrarily annotate the data with such attribute-value pairs, the users are often unwilling to perform this task: The task not only requires considerable effort but it also has unclear usefulness for subsequent searches in the future: who is going to use an arbitrary, undefined in a common schema, attribute type for future searches? But even when using a predetermined schema, when there are tens of potential fields that can be used, which of these fields are going to be useful for searching the database in the future?

Such difficulties results in very basic annotations, if any at all, those are often limited to simple keywords. Such simple annotations make the analysis and querying of the data cumbersome. Users are often limited to plain keyword searches, or have access to very basic annotation fields, such as "creation date" and "owner of document."

## 2. LITERATURE REVIEW

S.R. Jeffery, M.J. Franklin, and A.Y. Halevy[1]In this paper the writer have proposed a paper Pay-as-You-Go User Feedback for Data space Systems This system propose a system which is a line of work towards using more expressive queries that leverage annotations is the "pay-as – you – go " querying strategy in data spaces. In data spaces users provide data integration hints at querying time. But in this paper it is assumed that data sources already contain structured information and the problem is to match the query attributes with the source attribute.

K. Saleem, S. Luis, Y. Deng, S.-C. Chen, V. Hristidis, and T. Li [2] in this paper they proposed a paper "Towards a Business Continuity Information Network for Rapid Disaster Recovery In this paper they consider the Crisis Management and Disaster Recovery have gained immense importance in the wake of recent man and nature inflicted calamities. They proposed a solution or model for pre-disaster preparation and post disaster business continuity/rapid recovery. In case of

disaster need of rapid information retrieval and sharing increases. This paper proposed a disaster management model which works well at some extent but it is not considering the effective retrieval.

R.T. Clemen and R.L. Winkler [4] in this paper proposed a paper "Unanimity and Compromise among Probability Forecasters" In this paper they work on probabilities of particular uncertain event. This helps us to find out annotation and attributes.

G. Tsoumakas and I. Vlahavas [5] in this paper propose a paper Random K-Label sets: An Ensemble Method for Multilevel Classification. This paper proposes an ensemble method for multilevel classification. The Random k-label sets (RAKEL) algorithm constructs each member of the ensemble by considering a small random subset of labels and learning a single-label classifier for the prediction of each element in the power set of this subset. In this way, the proposed algorithm aims to take into account label correlations using single-label classifiers that are applied on subtasks with manageable number of labels and adequate number of examples per label. Using this we can take into account the correlation between tags for annotations. But in this collaborative annotation is missing.

S.-T. Wu, Y. Li, Y. Xu, B. Pham, and P. Chen [13] in this paper propose a paper "Automatic Pattern-Taxonomy Extraction for Web Mining," and "Deploying Approaches for Pattern Refinement in Text Mining," In these papers a technique of closed sequential patterns is used in text mining. It contains the concept of closed patterns in text mining. It improves the performance of text mining. Pattern taxonomy model is developed to improve the effectiveness. It uses closed patterns in text mining effectively. Term-based methods and pattern based methods is used to improve the performance of information filtering.

## 3. RELATED WORK

Existing work is done in this work format where the annotation scheme is being improved by CADS technique; here is the flow how the existing flow is working.
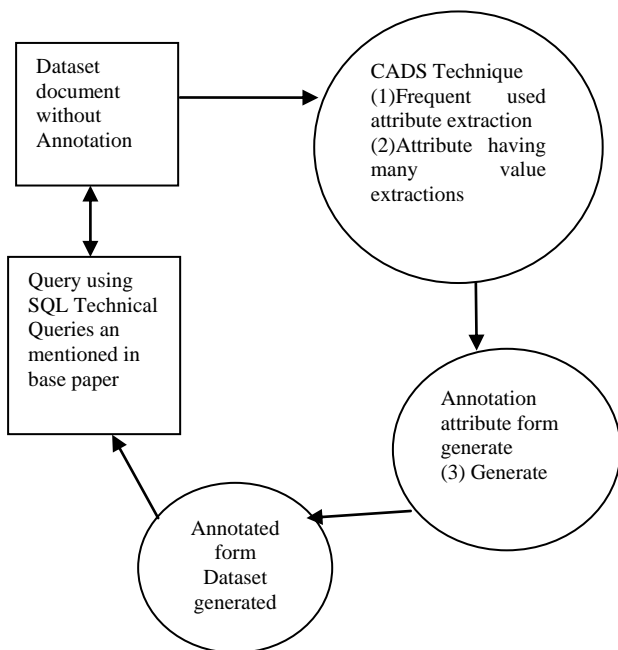


**Fig-Existing System Flow**

## 4. ISSUES WITH PREVIOUS TECHNIQUE

Problem with the existing work is while the efficient technique for annotation is available but still the data can only be access by the technical expertise but in the situation on demand where the possibility for the user to get the idea to access the structured data in need to avail for non-technical candidates who are often on demand requirement of these data of related to disaster or any sort of result which they often need to enquire on phone or via some short info.

The following issues have been arising with the existing technique:

- No efficient query algorithm is proposed yet.

- The available query technique is technical scope specific.

- Result can only be obtained based on the accurate query.

- No knowledge of data type without any prior study or knowledge about the accessed data on which user is working.

## 5. PROPOSED WORK

A Enhancement to the CADS where the efficient querying form also make generate where the querying to the structured dataset can be used by the normally people who might interested to access the data which is available for them, the data might be querying in the manner so that can be used by various research fellows and by different department of related data, we further can compare different querying technique by other approaches and experimentally can perform our efficient work over other work, because here we observed that number of work has been done for the structuring of dataset but still accessing them is technical task which is not for non-technical people, so here further we can extend the technique to work upon querying the dataset and also we can perform other querying technique and can compute precision, recall, query value, Computation time and computation cost.

We can definitely get good results in all the aspects as the computation time also will get less as less recourse need to provide while the data need to publish in public.

Here efficiency for every technique is often calculate using recall and precision and use to observe the results by different image technique while querying the dataset which we retrieved after annotation operation.

For query q,

A(q)=A set of data in dataset

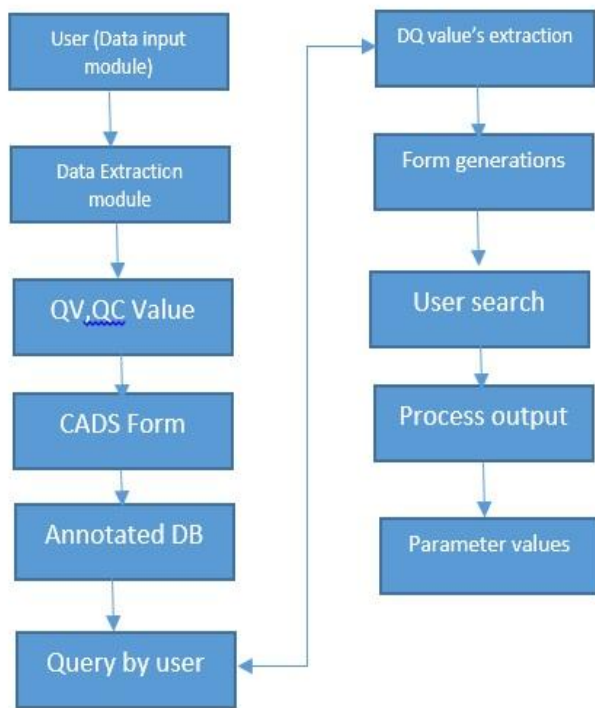B(q)=A set of data relevant to query

Precision = $A(q)\wedge B(q)/A(q)$

Recall=$A(q)\wedge B(q)/B(q)$

So precision is represents the ratio of the number of images relevant to the query q among retrieved data to the number of retrieved images & Recall is the ratio of the number of data relevant to the query among retrieved data to the number of data relevant to the query in a DB, so all the evaluation process of query the dataset are always done using these two calculations. These things we can perform with different retrieval technique on CADS.

- We are providing a form based searching technique where multiple users can create their own profile by following the login portal.

- User will get the same profile id like the admin.

- User can upload the dataset annotate the dataset along with that user can query from the all available datasets or from his own one.

- Users do not need to search again and again the last query because there is a "view" option through which user can view the last search.

- User can view all the uploaded datasets of admin along with the other users.

- Only admin will have the right of deletion of any categories or words.



**Flow chart**

**Algorithm Used**

**Input**: SDS, Models, Accuracy

**Output**: Query output & QCADS Form

**Steps:**
1. Retrieve next Qn from column 1-n.
2. Get the column and data value for each column $c_i$ Calculate QQV   QCV.
3. Calculate QCV.
4. Constant or Calculate T = avg (QCV). Where QCV is the maximum possible extraction of the all unseen values.
5. Sequentially process data
   i. T> value
   ii. R   put value
6. Apply QPP process to get model list and adaptive accuracy which is provided in input.
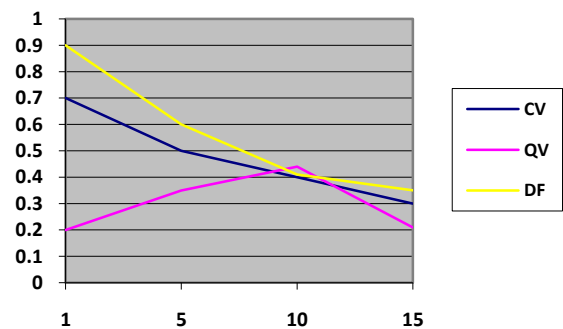7. Apply QCV
   a. Else

   b. Go to Step 1

**1.Evaluation of experiments:**
The experimental setup and the calculation of the efficiency for every technique is often calculate using recall and precision and use to observe the results by different image technique while querying the dataset which we retrieved after annotation operation.
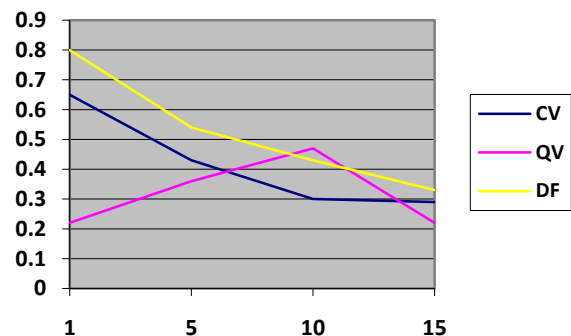
**Table 1: Dataset Statistics**

|  | DS1 | DS2 | DS3 |
|---|---|---|---|
| No. | Bhopal gas Tragedy | Historical places in India | Weather Forecast |
| Max(Size) | 26 KB | 9 KB | 56KB |
| Min(Size) | 3 KB | 2 KB | 9KB |
| Avg(Size) | 9 KB | 4.3 KB | 14KB |
| Annotation(Max) | 126 | 140 | 180 |
| Annotation (Min) | 16 | 12 | 8 |
| Annotation per Document(Avg) | 81 | 92 | 96 |

So precision is represents the ratio of the number of images relevant to the query q among retrieved images to the number of retrieved images & Recall is the ratio of the number of images relevant to the query among retrieved images to the number of images relevant to the query in a DB, so all the evaluation process of query the dataset are always done using these two calculations. These things we can perform with different retrieval technique on CADS.
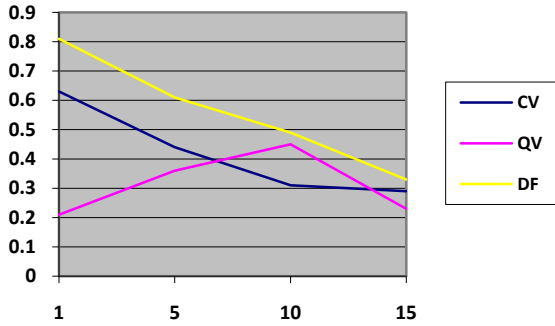


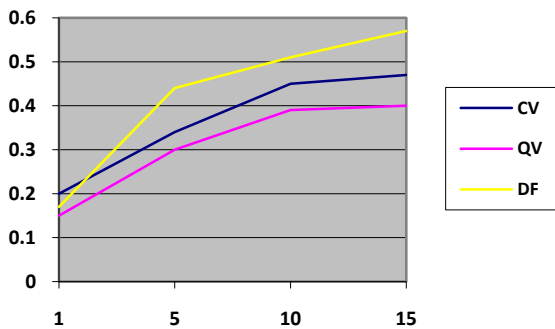Precision for Bhopal Gas Tragedy



Precision for Historical places India

Precision for News Blog



Recall for Bhopal Gas Tragedy



Recall for Historical Places India



Recall for News Blog

**Computation time**

Computation time is the length of time required to perform a computational process. Computation time can be represented as a sequence of time slots for performing

computation on the various available segments of the services. The computation time is proportional to the number of services.
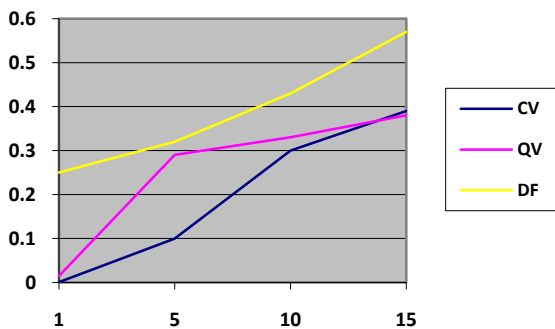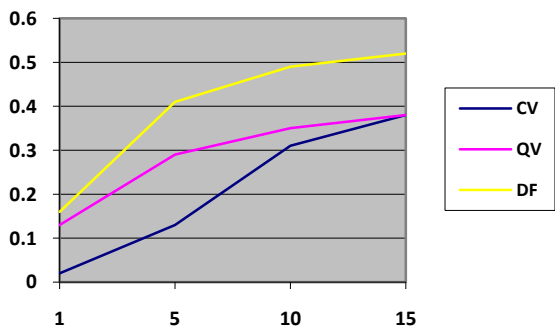
# 6. RESULT ANALYSIS

In this implementation, we included two datasets of weather forecasting and Bhopal gas tragedy. Weather forecasting has 44 files and Bhopal gas tragedy has 17 files. On basis of these datasets annotation performed. Annotation are dynamic in nature as if new document is added, is matching keywords are available, it will automatically annotate it and if need to edit, it can. In the graph representation, annotation, keywords, files, counts, top keywords, top values and top used keywords are used to illustrate.
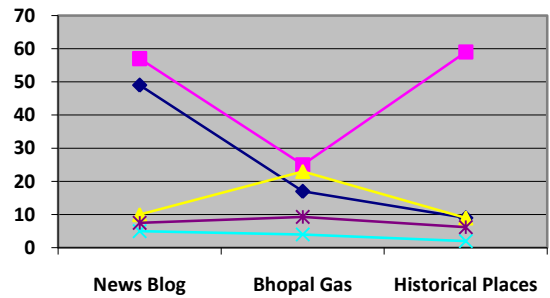


Annotation, Keywords & Files



Result time (in seconds)



Top count values

Top count keywords



## 7. CONCLUSION

We have discussed the existing technique which is using in the annotation world in order to annotate the data and querying the data according to the user requirements, we also saw the schemes are used in the case of annotation access in the data annotation and access world. As per the discussion and the work. we have got to know about the work which is already done in the field of annotation and accessing the dataset which is keep structured with the help of different annotation based algorithm, we have also discussed CADS algorithm which is proposed in our base paper and also we have mentioned the limitation of the existing & available algorithm, so here upon discussion we are proposing new algorithm which is efficient and going to work in the field of accessing annotated dataset.

## 8. REFERENCES

[1] [1] S.R. Jeffery, M.J. Franklin, and A.Y. Halevy: proposed a paper "Pay-as-You-Go User Feedback for Data space Systems,"

[2] K. Saleem, S. Luis, Y. Deng, S.-C. Chen, V. Hristidis, and T. Li: proposed a paper "Towards a Business Continuity InformationNetwork for Rapid Disaster Recovery.

[3] J. M. Ponte and W.B. Croft: proposed a paper "A Language Modeling Approach to Information Retrieval".

[4] R. T. Clemen and R.L. Winkler: proposed a paper "Unanimity and Compromise among Probability Forecasters.

[5] G. Tsoumakas and I. Vlahavas: propose a paper "Random Label sets: An Ensemble Method for Multilevel Classification.

[6] P. Heymann, D. Ramage, and H. Garcia-Molina: proposed a paper "Social Tag Prediction".

[7] Y. Song, Z. Zhuang, H. Li, Q. Zhao, J. Li, W.-C. Lee, and C.L. Giles: proposed a paper "Real-Time Automatic TagRecommendation".

[8] D. Eck, P. Lamere, T. Bertin-Mahieux, and S. Green: proposed a paper "Automatic Generation of Social Tags for MusicRecommendation.

[9] B. Sigurbjornsson and R. van Zwol: proposed a paper "Flickr Tag Recommendation Based on Collective Knowledge".

[10] B. Russell, A. Torralba, K. Murphy, and W. Freeman: propose a paper "Label Me: A Database and Web-Based Tool for ImageAnnotation".

[11] M. Franklin, A. Halevy, and D. Maier: propose a paper "From Databases to Data spaces: A New Abstraction for InformationManagement ".

[12] J. Madhavan et al: proposed a paper "Web-Scale Data Integration: You Can Only Afford to Pay as You Go".

[13] S.-T. Wu, Y. Li, and Y. Xu, "Deploying Approaches for Pattern Refinement in Text Mining," Proc. IEEE Sixth Int'l Conf. DataMining (ICDM '06), pp. 1157-1161, 2006.

[14] D. Liu, X.-S. Hua, L. Yang, M. Wang, and H.-J. Zhang, "Tag Ranking," Proc. 18th Int'l Conf. World Wide Web (WWW), 2009.

[15] D. Yin, Z. Xue, L. Hong, and B.D. Davison, "A Probabilistic Model for Personalized Tag Prediction," Proc. ACM SIGKDD Int'l Conf. Knowledge Discovery Data Mining, 2010.

[16] K. Chen, H. Chen, N. Conway, J.M. Heller stein, and T.S. Parikh, "Usher: Improving Data Quality with Dynamic Forms," Proc. IEEE 26th Int'l Conf. Data Eng. (ICDE), 2010.

[17] M. Jayapandian and H.V. Jagadish, "Automated Creation of a Forms-Based Database Query Interface," Proc.VLDB Endowment, vol. 1, pp. 695-709, Aug 2008.