# Deep web interface for Fine-Grained Knowledge Sharing in Collaborative Environment

Andrea.L[1] , S.Sasikumar[2]

[1]*PG.Scholar,*
*Department of Computer Science and Engineering*
*Saveetha Engineering college*
*Tamilnadu , India*

[2]*Professor ,*
*Department of Computer Science and Engineering*
*Saveetha Engineering college*
*Tamilnadu , India*

***Abstract*: As profound web grows at a very fast pace, there has been an increased interest in techniques that help efficiently to locate deep-web interfaces. Due to the huge volume of web resources and its dynamic nature of deep web, achieving wide coverage and high efficiency is a tedious issue. In proposed system the multi-keyword search concept is used , where the system will be providing all possible relevant links. The search is based on each word of the input data. The query which is submitted to the application will be preprocessed based on ontology word net tool, after pre-processing only root words will be taken and Synonym, Hypernym and Hyponym of the root words will be listed to the user. The links relevant to the search functionality are being categorized under various groups. The links obtained during the search process are bookmarked either locally or globally. The bookmarks saved locally are browser independent. The universal bookmarks are available to the users and are viewed based on search query. The main goal is to find out the knowledgeable member and to connect the member to the search user to obtain more information.**

***Keywords*:Deepweb, pre-processing ,Synonym , Hypernym Hyponym , Bookmark.**

## I INTRODUCTION

The deep web alludes to the substance lie behind searchable web interfaces that can't be recorded via search engines. To cluster the content archive over the page in view of the client input key term, to upgrade profound web look and to overcome gathering of inconsequential records into the same group is an extremely challenging procedure. The proposed framework intends to Web-clients to find the best method for their pursuit needs, bringing about speedier and more precise query items.

A feasible profound web gathering structure named SmartCrawler is being for fulfilling both wide extension and high viability for a focused crawler. Smart crawler is isolated into two stages: site discovering and in-site exploring [1]. The webpage finding stage achieves wide extent of destinations for a connected with crawler, and the in-webpage examining stage can adequately perform looks for web outlines within a site.

A novel two-stage structure is being utilized to address the issue of searching for hid web resources. Our site discovering technique uses an opposite looking strategy and incremental two-level site sorting out framework for revealing applicable destinations, finishing more data sources. An adaptable learning computation that performs online part determination and utilizes these segments to normally fabricate join rankers. In the site finding stage, high huge destinations are sorted out and the inching is revolved around a subject using the substance of the root page of districts, fulfilling more exact results. Amid the insite examining stage, appropriate associations are sorted out for snappy in-site looking.

To influence the vast volume data covered in deep web, past work has proposed various procedures and devices, including profound web understanding what's more, reconciliation , hiddenweb crawlers , and deep web samplers .For all these methodologies, the capacity to creep deep web is a key test. Olston and Najork efficiently exhibit that slithering profound web has three stages: finding profound web content

sources, selecting significant sources and extricating hidden content . Taking after their announcement, we examine the two stages firmly identified with our work as underneath. Finding profound web content source. Generic crawlers are for the most part produced for describing deep web and registry development of deep web assets, that don't restrain seek on a particular theme, yet endeavour to get all searchable shapes . The Database Crawler in the MetaQuerier is intended for consequently finding question interfaces. Database Crawler first finds root pages by an IP-based inspecting, and after that performs shallow slithering to creep pages inside of a web server beginning from a given root page .The FFC contains three classifiers: a page classifier that scores the importance of recovered pages with a particular subject, a connection classifier that organizes the connections that might prompt pages with searchable structures, and a structure classifier that sift through non-searchable structures.

The enhancement FFC is considered with an versatile connection learner and programmed highlight choice. SourceRank , evaluates the pertinence of profound web sources amid recovery. Taking into account an assention diagram, SourceRank computes the stationary visit likelihood of an irregular stroll to rank results. Not quite the same as the creeping systems and instruments said above, SmartCrawler is a space particular crawler for finding significant.

## II RELATED WORK

As profound web develops at a quick pace, there has been expanded enthusiasm for procedures that help effectively find profound web interfaces. Due to the huge volume of web resources and its dynamic way of profound web, achieving wide scope and high effectiveness is a challenging issue. The expanding number of ontology's of the Semantic Web postures new difficulties for mapping[6].

A key objective of the Semantic Web is to move social connection designs from a maker driven worldview to a purchaser driven one [19]. It conceptualizes the customization of administration based business forms utilizing the current information of Web administrations and business forms. Extensible Mark-up Language (XML) mark up dialect Web Ontology Language-Business Process Customization (OWL-BPC), taking into account the accepted semantic mark-up dialect for Web-based data [Web Ontology Language (OWL)].

While trying to find the profound web databases, in light of the fact that they are not enlisted with any search engine. To address this issue, past work has proposed two sorts of crawlers namely generic crawlers and focused crawlers. Generic crawler and Focused crawler.

Generic crawler cannot concentrate on a particular subject and is time consuming .Focused crawlers concentrates on particular subject but does not follow any orderly standards or ranking procedures. The highly ranked search are not being viewed in an orderly manner. These crawlers are a major

drawback for providing search data to the user. More over during the search functionality irrelevant data is produced, which should be taken care and should be avoided.

## III ARCHITECTURE OF DEEP WEB INTERFACE

The Architecture diagram of the deep web interface is shown below. The user has to register in the web application. After successful registration the user is directed to the Search page where the user can provide the input data to be search.
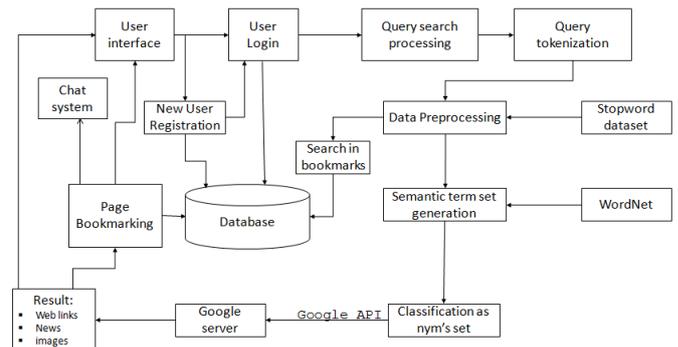


Fig 1 : Working of deep web interface

Data pre-processing functionality enables the removal of the stop words . The relationship for the given input data is being found using the ontology Word net tool .Search functionality is done word by word and not based on keyword based search. Each word is given equal priority. Poster stemming procedure alludes to creating a lookup table which contains relations between root forms and inflected forms. Ontology clustering uses the nym's classification to describe different classes of words , the relationship between words and categorize the result obtained under three groups namely Synonym , Hyponym and Hypernym .This functionality enables non-technical user to understand complicated concepts by knowing their common meaning. According to the user's needs the bookmarking takes place either locally or globally. The links obtained after search are bookmarked either locally or globally. The bookmarks are browser independent. The global bookmark enables to find a knowledgeable person and connects the knowledgeable person to the search user to obtain more information.

## IV SYSTEM IMPLEMENTATION

### A. User Interface

An usual application with User registration page enables the new user account creation. Mail-id and mobile number are the unique entities being validated during the new user registration process. Authentication of the existing users

**ISSN: 2278 – 1323**

*International Journal of Advanced Research in Computer Engineering & Technology (IJARCET)*
*Volume 5, Issue 4, April 2016*

enhances security. The user needs to register to prevail the search functionality usage.

After the successful authentication , the web user is directed to the search space page. This is the environment for user to search the content from the web server. This Search Space act as the interface between the user and the web servers. The input data to be searched is being entered by the user in the search space provided.

The user provides the input text to be searched. Based on the input text the search functionality begins where pre-processing takes place word by word.

### B. Data Pre-processing

User interface functionality is followed by the data pre-processing. Stop words are the words which are removed out prior to, or after processing of natural language data (text). It is only controlled by human input and not automated. The following are some of the most common stop words , such as "the, is, at, which ,who ,why, when ,above, below  and on". The stop word removal functionality enhances the speed of the search result.

Stemmer algorithm  used in this process, employs a lookup table which gathers the relationship between root forms and inflected forms. To stem a word, the lookup table is queried to find out a matching word. If a matching word is found, the respective root word is returned. Eg: A stemming algorithm reduces the words "sailing", "sailed" , "sail", and "sailor" to the root word , "sail".

### *Steps involved in Poster Stemming algorithm*

Step 1 : Get  rid of the plurals and –ed or –ing  suffixes.
Step 2 : Turns terminal y to I when there is  another vowel in the stem.
Step 3 : Maps double suffixes to single ones: -ization, -ational, etc.
Step 4 : Deals with suffixes, -full , -ness, etc.
Step 5 : Takes off-ant, -ence, etc.
Step 6:  Removes a final –e.

### C. Ontology Clustering

Data pre-processing process is followed by the ontology clustering. Ontology deals with the relationship between the word existence. Words ending in "nym's" are used to describe different contexts of words and  the relationship between words.
- Hypernym: A word that has a general meaning than another.
- Hyponym: A word that has a specific meaning than another.
- Synonym : One of two (or more) words that have the same (or very similar)

Ontology word net tool is being used for finding the relationship between the words. Based on the relationship the word clustering takes place.

The Artificial-Intelligence literature contains various many definitions for ontology (Wordnet) in finding the relationship between the words .It includes machine-interpretable definitions of basic concepts of the domain and the relationship among  them. The obtained  featured results produced by the sentence-based, document-based, corpus-based, have higher quality than those produced by a single-term analysis similarity.

### D. Multi-term search

In the existing system search functionality was done based on keywords. Only the  search links related to the keywords were only produces. In Deep web interface system search results are produced for each and every word after the pre-processing activity. The users provide the multi-term input to be searched. Multi-term search provides the result deeply from the search engines and its search the terms randomly till last key term in that multi-term list.

### E. Clustering and Bookmarking

From the multi-term search result we cluster the more relevant content based on the relationship user input term. And we classify the cluster and give the final output like most relevant content comes first and out comes next output screen .The result obtained are features as news, web links and images. Search functionality takes place more deeply as the least search information are also produced. Usually while searching only the links which are viewed more frequently and produced as results and the least viewed links are not specified. For eg. If we search for a poet or guitarist  who has the same name of a sport's person , the result obtained will be of the sport star. The poet or the guitarist will not be listed out as they are least viewed .This functionality is being overcome in the deep web search , as all the related search are listed out based on the ranking functionality.

The result obtained is bookmarked either locally or globally according to the user's need. The links bookmarked are browser independent. The universal bookmark are made available to the registered users according to the search made. The universal bookmark are ranked according to the user's review and preference. The main aim of the functionality is to connect the knowledgeable person to the search user to obtain more information in a specific domain.

## V RESULT ANALYSIS

The existing system is based on the keyword search functionality. Only the keywords are picked from the input data given by the user and the relevant links are produced. Due to this feature only surface search takes place and the

1084

deep search of gaining more information is suspended. The proposed system provides multi-term search functionality where search functionality takes place word by word. Once the user provides the input data to be searched the data pre-processing takes place. After the removal of the stopwords only the root words are retained for the search functionality
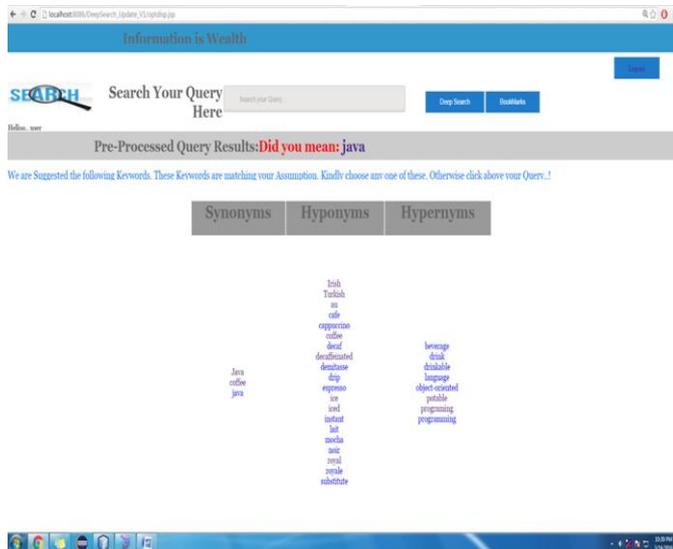


Fig 2 : Classification using nym's group

The root words are the processed for ontology clustering where the relationship between the words are obtained. Based on the classification the result are categorized under three groups namely Hypermym , Hyponym and Synonym. Each link provides data under web links, news and images. The search links are displayed in ranking order. Based on the user's preference the links are bookmarked locally or globally. The locally saved bookmarks are browser independent. The universal bookmarks are available to the registered users. The universal bookmarks are helpful in finding a knowledgeable person to get more detailed information for a particular field or domain.

## VI CONCLUSION AND FUTURE WORK

Web mining plays an important role in the extraction of information . It is an area full of challenges and many research problems are yet to be identified. Various techniques and methods are used in the extraction of information from Deep web. The proposed system make use of the multi-term search functionality with ontology classification and groups the result obtained under nym's groups . By doing so more enhanced and huge variety of data are obtained during search. Bookmarking plays a vital role as they are browser Independent . The bookmarks at present are browser independent. As the browser crashes all the information stored are also lost. Universal bookmarks are made available to the

registered users that paves the way to find a knowledgeable person and to connect the person to the search user.

Future work it can be implemented by providing a direct video chat or audio chat to obtain more detailed information.

## REFERENCES

[1] Feng Zhao , Jingyu Zhou , Chang Nie , Heqing Huang, Hai Jin (2015) ," A Two-stage Crawler For Efficiently Harvesting Deep - Web Interfaces".

[2] Vinay Kancherla (2014) , " A Smart Web Crawler For a Concept Based Semantic Search Engine".

[3] Thirumaran.M , Dhavachelvan.P , Shanmugapriya , Aishwarya.D (2013) ,"Ontology based Dynamic Customization for Composite Web Services".

[4] Michael Rys , Ka Fai Yau , " Data Extraction from Dynamic Web Sites : Combining Crawling and Extraction".

[5] C´assia Trojahn , Paulo Quaresma, Renata Vieira ," A Framework for Multilingual Ontology Mapping ".

[6] M. Nagy and M. Vargas-Vera(2011), "Multiagent ontology mapping framework for the semantic web," IEEE Trans. Syst., Man, Cybern. A, Syst., Humans, vol. 41, no. 4, pp. 693–704.

[7] Antonio Sanfilippo, Stephen Tratz, Michelle Gregory, Alan Chappell, Paul Whitney, Christian Posse, Patrick Paulson, Bob Baddeley, Ryan Hohimer, Amanda White , "Ontological Annotation with WordNet".

[8] Mingming Li1, Chunlin Li1, Chao Wu1 and Youlong Luo2 ," A Focused Crawler URL Analysis Algorithm based on Semantic Content and Link Clustering in Cloud Environment ".

[9] Bo Fu1, Rob Brennan1 and Declan O'Sullivan," Multilingual Ontology Mapping: Challenges and a Proposed Framework".

[10] Christina Brandt, Thorsten Joachims, Yisong Yue, Jacob Bank," Dynamic Ranked Retrieval".

[11] Andr´e Bergholz and Boris Childlovskii (2000),"Crawling for domain specific hidden web ". In Proceedings of the 27th International Conference on Very Large Data Bases, pages 129–138.

[12] Balakrishnan Raju and Kambhampati Subbarao. Sourcerank(2011),"Relevance and trust assessment for deep web sources based on inter-source agreement". In Proceedings of the 20th international conference on World Wide Web, pages 227–236.

[13] Balakrishnan Raju, Kambhampati Subbarao, and Jha Manishkumar(2013)," Assessing relevance and trust of the deep web sources and results based on inter-source agreement". ACM Transactions on the Web, 7(2):Article 11, 1–32.

[14] Brightplanet's searchable database directory(2001). http://www.completeplanet.com/, 2001.

[15] Cheng Sheng, Nan Zhang, Yufei Tao, and Xin Jin(2012),” Optimal algorithms for crawling a hidden

[16] Christopher Olston and Marc Najork(2010),”Web Crawling”. Foundations and TrendsR in Information Retrieval Vol. 4, No. 3, 175–246c 2010 C.

[17] Denis Shestakov and Tapio Salakoski(2007),” On estimating the scale of national deep web”. In Database and Expert Systems Applications, pages 780–789. Springer.

[18] Q. Liang, X. Wu, E. K. Park, T. M. Khoshgoftaar, and C. H. Chi (2011), "Ontology-based business process customization for composite web services," IEEE Trans. Syst., Man, Cybern. A, Syst., Humans, vol. 41, no. 4, pp. 717–729.

[19] M. Cai, W. Y. Zhang, and K. Zhang (2011), "ManuHub: A semantic web system for ontology-based service management in distributed manufacturing environments," IEEE Trans. Syst., Man, Cybern. A, Syst., Humans, vol. 41, no. 3, pp. 574–582.

[20] Jayant Madhavan, Shawn R. Jeffery, Shirley Cohen, Xin Dong David Ko, Cong Yu, and Alon Halevy(2007),” Web-scale data integration”.In Proceedings of CIDR, pages 342–350.

database in the web”. Proceedings of the VLDB Endowment, 5(11):1112–1123.