



ISSN(Online) : 2320-9801  
ISSN (Print) : 2320-9798

## International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Vol. 4, Issue 4, April 2016

# Smart Web Crawler to Harvest the Invisible Web World

Nimisha Jain<sup>1</sup>, Pragya Sharma<sup>2</sup>, Saloni Poddar<sup>3</sup>, Shikha Rani<sup>4</sup>

B.E. Student, Dept. of Computer Engineering, MIT College of Engineering, Pune, India<sup>1,2,3,4</sup>

**ABSTRACT:** In a growing world of internet, thousands of web pages are added daily. But only approx. 0.03 percent fraction of web pages are retrieved by all the search engines. The remaining pages are deep websites. Deep web is that part of web which is hidden and unrecognizable by the existing search engines. It has been a longstanding challenge for the existing useful crawlers and web search engines to harvest this ample volume of data. This paper surveys on different methods of crawling deep-web. Density and rapidly changing nature of deep-web has posed a big hurdle in front of researchers. To overcome this issue, we propose a dual-layer framework, namely *Smart Web Crawler*, for efficiently harvesting deep web interfaces. Smart crawler consists of two layers: *Site-discovery* and *In-depth Crawling*. *Site-discovery* layer finds the sparsely located deep websites from given known *parent sites* using *Reverse Searching* and focused crawling. The *In-depth Crawling* layer makes use of *Smart Learning* and *Prioritizing* to crawl hyperlinks within a site to ensure wider coverage of web directories.

**KEYWORDS:** Deep Website, HTML Form, Reverse Searching, Prioritizing, Smart Learning, Categorizing, Pre-querying, Post-querying

### I. INTRODUCTION

The basic motive behind web crawler is to retrieve web pages and add them to a local repository. Such repositories are used with web search engines. Basically it starts from a parentsite and then uses hyperlinks and sub-links contained in it to reach other pages. It records every hypertext link in every page they index crawling. It keeps on repeating until it reaches a predefined value.

Today Internet can be considered as a vast pool of information from which any data can be found. But still there is so much data which cannot be retrieved by any search engine. Such data is called deep web pages. The deep web refers to the contents that lie behind HTML forms and cannot be identified by normal crawlers. This data is buried deep down on dynamic web pages and hence cannot be found by any search engines. They include dynamic web pages, blocked sites, unlinked sites, private site content and limited-access networks. The information available on deep web is approximately 400-550 times larger than the surface web. More than 200,000 deep websites presently exist [1].

Current-day crawlers can retrieve only websites which are indexed, i.e., the webpages that can be reached by following hypertext links, ignoring webpages that need authentication or registration. So, they ignore a large amount of high quality content hidden behind search forms, in large searchable electronic databases. Deep web contains huge amount of valuable information which is highly scattered. This information may be required to build index in a given domain by entities such as Infomine, Clusty. But these entities do not have access to the licensed web indices of standard search engines. So, there is a need of a proficient crawler which can mine deep web interfaces quickly and accurately. Traditional search engines cannot retrieve content in the deep Web as they are created dynamically. The Deep Web offers a certain level of obscurity of identity that makes it more vulnerable to be used for illegal purposes. The various activities that take place in it, such as illegal drug selling, child pornography etc., show what people would do if their identities would be hidden. Instead of just considering the search results from individual we can train a crawler to learn from features collected from path of connected links.



# International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Vol. 4, Issue 4, April 2016

## II. RELATED WORK

In [2, 3], author outlined that crawler is a program which is used to download or store web pages by indexing them. Mainly it starts with maintaining a set of URLs to be retrieved and prioritized in a queue. From this queue, the crawler gets a URL to be downloaded. Each URL is called a *parent*. It checks if the URL or page is allowed to be downloaded and also reads the header of the page to check if any exclusion instructions were provided. It extracts all hyperlinks from the page and add them to the queue used for storing the URLs to be downloaded. This queue is termed as *crawl frontier*. Using a set of rules and policies the URLs in the frontier are visited individually. Then it downloads different pages from the internet by the parser and generator. These pages are stored in the database system of the search engine. The URLs are then placed in the queue and later scheduled by the scheduler and can be accessed one by one by the search engine one by one whenever required. Various algorithms are used for finding the relevant links. These relevant links are prioritized according to the page ranks. Page ranks are assigned according to the relevance of the page. The algorithms used are Breadth first algorithm, Naïve Best first algorithm, Page ranking algorithm, Fish search algorithm etc. [4] There are various types of crawlers currently available for various operations such as *Internet archive Crawler*, *Google Crawler*, *Mercator Web Crawler*, *Web Fountain crawler*, *IRLbot Web crawler*, *Hidden Web Crawler* etc.[5].

In [6], the goal of a focused crawler is to select topic-specific pages that are relevant. Instead of crawling and indexing entire web, it focuses on traversing the links within the predefined boundary and thus avoids crawling irrelevant regions. This leads to efficient usage of hardware and network resources. Use of multiple focused crawlers result in entire coverage of web at a faster rate as each crawler is dedicated to a particular domain and later guided by a classifier and a distiller.

In [7], the issues related to dual essential tasks are demonstrated: First, for *exploration*, our *macro* study surveys the deep web at large: It studies web databases in a broader sense, without giving much details about them, such as: What is its scale? How many databases are there? Where to find entrance to them? How many are structured databases? What is the coverage of deep-web directories? Second, for *integration*, our *micro* study surveys domain characteristics in much detail, such as: How hidden are web sources? How do search engines cover their data? How complex are their directories? What is the category distribution of sources?

In [8], author addresses the problem of designing a crawler which can be used to crawl the concealed web world. To achieve this, a simple operational model which describes the steps that a crawler must take to process and submit HTML forms. Based upon the working of this model, a prototype crawler HiWE (Hidden Web Exposer) is introduced. It ensures feasibility and effectiveness of our form processing and matching techniques. This model makes use of a new Layout-based Information Extraction Technique (LITE) to automatically extract pertinent information from search forms and response pages.

In [9], it is discussed that as the scale and heterogeneity is increases, the traditional data integration techniques are no longer valid. Thus, we propose a new data integration architecture, the PAYGO architecture which is inspired by the concept of dataspace [10]. It comprises of two steps based on pay-as-you-go data management: First, we focus our work on searching the deep web. This approach leaves the data at the sources and directs queries to appropriate web pages, attempts to add data from the deep web to surface web and also tells about the first study of the deep web which is based on a commercial index. Second, we consider Google Base and define how a model can be used to improvise user's search experience.

In [11], the problem of distributed information retrieval and its corresponding solutions are discussed. We describe a technique for identifying search forms, which builds the base for a next-generation distributed search problem. It provides a unified search experience to the users using available search interfaces and displays the results of the queries into an integrated list. Automatic feature generation technique is used to generate candidate forms and C4.5 decision trees to classify them. This improves the accuracy as well as the precision.

This [12] resolves the major challenges in building Deep Web integration system using VisQI (VISual Query interface Integration system). VisQI is a framework-like reusable architecture that compares generated data structures against gold standard. It transforms the HTML query interfaces into ecclesiastic structure, classifies them into different domains and matches the elements of various interfaces.

In [13], the author talks about a new crawling strategy to automatically discover concealed web databases to achieve a balance between the two opposing requirements of this problem: that is performing a broad search while avoiding the need to crawl a large number of insignificant pages. It performs focused crawling using appropriate stopping criteria

# International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Vol. 4, Issue 4, April 2016

crawler to avoid excessive speculative work in a single website. For ranking we use the algorithm based upon the function of the number of pages visited by the user. We use backward crawling facilities in order to approximate web graphs.

### III. PROPOSED WORK

In this paper we are proposing a *Smart Web Crawler*, which is a dual-layered architecture to harvest the deep and hidden web interfaces. To rank unexplored links on priority basis, *Site Prioritizer* is used in the first layer. When the crawler finds a new deep website, its URL is added into the *Database*. A *Smart Learning* technique is used where the crawler learns from the URL path leading to relevant forms, thus enhancing the performance of the *Prioritizer*.

To ensure wide coverage of deep web interfaces and efficient crawling, our crawler consists of two layers: *Site-discovery* and *In-depth Crawling*. *Site-discovery* layer finds the sparsely located deep websites from given known *parent sites* using *Reverse Searching* and focused crawling. The *In-depth Crawling* layer makes use of *Smart Learning* to crawl hyperlinks within a site to ensure wider coverage of web directories. To improve accuracy of *HTML form categorizer*, two approaches are used: *Pre-querying* and *Post-querying*.

#### A. Design:

For efficiently harvesting deep web data directories, *Smart Web Crawler* is designed for deep websites. It has a dual-layered architecture. The first layer is *Site-discovery* and the second layer is *In-depth Crawling* as shown in figure. *Site-discovery* layer finds the most appropriate and relevant sites for a given query by user, and in the second layer *In-depth Crawling* layer covers hidden HTML forms from the site.

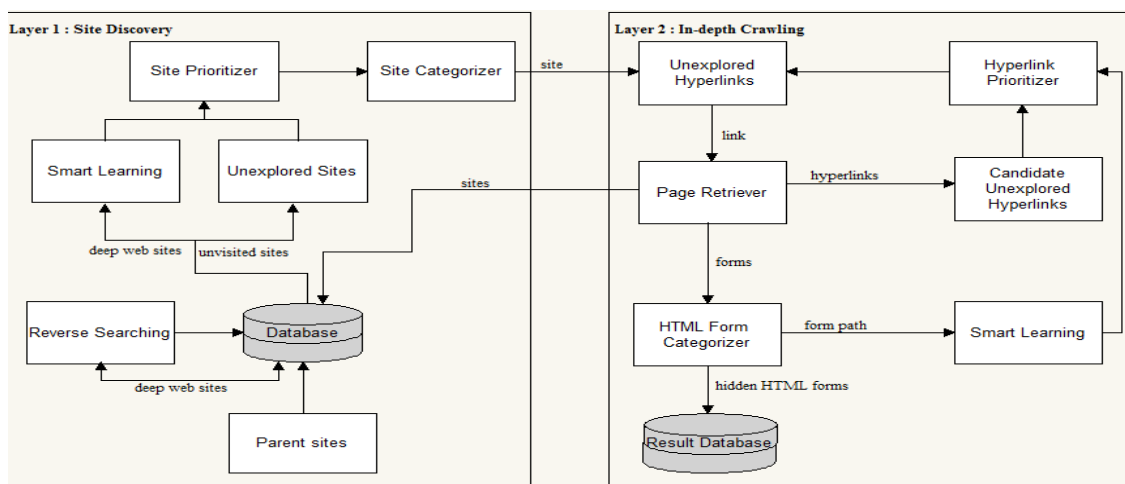


Figure 1: Proposed Architecture

#### 1) Site Discovery

Initially, the *Site-discovery* layer begins with a *parent URLs* of sites in database. *Parent URLs* are primary URLs given to *Smart Web Crawler* to start crawling. It starts by following URLs from selected *parent sites* to search available web pages and domains. *Reverse Searching* is performed when the number of *Unexplored URLs* becomes less than the predefined value. In *Reverse Searching*, it performs the crawling of the known deep websites to find highly ranked center pages which point to several other hyperlinks/domains and adds them into *Database* simultaneously. The *Unexplored URLs* are ranked by the *Prioritizer* and visited URLs are added into the retrieved site list. The *Prioritizer* associates a value with each unexplored site based upon its content with respect to the already discovered deep web interfaces. Meanwhile, the *Prioritizer* is improved during crawling by a *Smart Learner* which continuously learns from features of already searched deep websites. To improve the accuracy of the search results, *Categorizer* classifies URLs into useful or useless for a given topic with respect to the homepage content.



# International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Vol. 4, Issue 4, April 2016

- *Parent Sites*: These are the URLs of deep websites given as input to start the crawling process.
- *Unexplored Sites*: It fetches homepage URLs from the *Database*, which are then ranked by *Prioritizer* according to the relevance. It also prioritizes the sub-links of visited web pages.
- *Reverse Searching*: When the number of unexplored sites decreases to a predefined value, *Reverse Searching* is triggered. The idea of *Reverse Searching* is to find center pages of unexplored websites starting from the last extracted hyperlink or sub-link to reach the initial *Parent Site* by crawling backwards. It ensures that there is enough number of unexplored URLs in the database.
- *Smart Learning*: The performance of the *Prioritizer* is improvised during the crawling process by using *Smart Learning* technique. It gradually learns from features of deep websites found and from the URL path leading to relevant hyperlinks.
- *Prioritizer*: The *Prioritizer* allocates a value to each and every unexplored site according to relevance in comparison with already discovered deep web sites.
- *Site Categorizer*: The crawler maintains a priority queue for all the sub-links and hyperlinks which are classified by *Categorizer* according to the relevance.

## 2) In-depth Crawling

After the most relevant site is found in the previous layer, this layer performs efficient *In-depth Crawling* for extracting hidden HTML forms. Hyperlinks pointed by the site are stored in the *Unexplored Hyperlink* section. And thus corresponding pages are fetched and then classified by *HTML form Categorizer* to retrieve hidden HTML forms [14, 15, 16, 17 and 18].

Additionally, the links in extracted centre pages from the previous layer are fetched into *Candidate unexplored URLs*. Here, they are prioritized by the *hyperlink Prioritizer*. Whenever a new site is discovered, it is first added into the *Database* and then ranked according to the relevance. The *Prioritizer* is improvised by using *Smart Learner* which gradually learns from the URL path leading to relevant websites.

To improve accuracy of *HTML form Categorizer*, *Pre-querying* and *Post-querying* approaches for classifying deep-web forms are combined. *Pre-querying* basically triggers before the query executes and is fired once while you try to query. With the help of this *Pre-querying* section, 'where' clause part can be changed dynamically. This query is fired only once. Whereas, *Post-querying* is triggered after the query executes gets fired for every row fetched. That is, if query fetches 10 rows then *Post-querying* will be fired for 10 times. Hence it is fired almost every time at the time of execution successfully during crawling.

- *Unexplored Hyperlinks*: Hyperlinks pointed by the most relevant sites fetched from previous layer are stored in this section and subsequently classified by *HTML Form Categorizer* to give hidden HTML forms.
- *Hyperlink Prioritizer*: It prioritizes the hyperlinks based upon the relevance of hidden HTML forms. It assigns relevance score to each hyperlink depending upon similarity with parent link.
- *Page Retriever*: It fetches the centre pages of websites. These pages are further classified to give relevant hidden HTML forms.
- *Candidate Unexplored Links*: It contains the *unexplored hyperlinks* of the *parent site*. These hyperlinks are prioritized and classified by the *Hyperlink Prioritizer*.
- *HTML Form Categorizer*: It ranks the *candidate unexplored hyperlinks* to give relevant hidden HTML forms and accordingly assigns a priority value. It filters out irrelevant or off-topic links.
- *Smart Learning*: It keeps on learning and adapting from the URL path leading to relevant HTML forms.
- *Result Database*: It is collection of all the relevant websites received by *HTML form Categorizer* at the end of whole process of crawling.

*Smart Crawler* implements a unique strategy to retrieve relevant hidden HTML forms by using various classifiers. It consists of two classifiers, a hidden HTML form classifier (HFC) and a field-specific form classifier (FFC). Hidden HTML form classifier is used to filter out irrelevant forms by matching their features with the standard layout. It is domain-independent. Field-specific form classifier fetches relevant forms using focused crawling.



# International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Vol. 4, Issue 4, April 2016

## IV. RESULTS

### A. Analysis of Proposed Algorithm:

In this subsection, we are analyzing the proposed algorithm in terms of time complexity. As quoted by Mathematician Alan Turing, “It is convenient to have a measure of the amount of work involved in a computing process, even though it be a very crude one. We may count up the number of times that various elementary operations are applied in the whole process”, we are performing analysis to determine the cost of each basic operation and hence develop a realistic model for the input. The total running time of any algorithm can be calculated as follows:

$$\sum \text{frequency}(a) * \text{cost}(a)$$

Where a= Operation

Following figure shows our prototype time-complexity analysis:

1. Begin	
2. Pick a site from the list of parent sites	$C_1$
3. Extract the page source	$d$
4. Parse the page source	$O(n)$
5. Get hyperlinks from the page source	$O(n)$
6. While(sites in hyperlink list is not empty) Loop is executed 'p' times	$O(m*p)$
7. If keyword is found in hyperlink, append it into result list	$O(m)$
8. Else discard it	$C_2$
9. Return hyperlink	$C_3$

Table 1: Prototype time-complexity analysis

From the above analysis, we have found time-complexity of our crawler as  $[2O(n)+O(m*p)+d]$ . Now that we know 'm\*n' is always  $\ll n$ , the above complexity can be re-written as  $[2O(n) + d]$ .

### B. Comparison Analysis of Different Crawlers:

In this subsection, we are comparing the strengths and weaknesses of the proposed crawler with some of the other crawlers. The purpose behind this comparison is to analyze and confront their performances. This gives the clear idea about feasibility and the efficiency of the crawler with respect to other crawlers.

Paper	Strength	Weakness
Raghavan et.al.[21]	The feasibility of HiWE algorithm and effectiveness of form processing and matching techniques.	1) Inability to recognize and respond to simple dependencies between form elements. 2) Lack of support for partially filling out forms
Liddle et.al.[22]	Automatically retrieves the data behind a given Web form.	1) Determining how best to partition the search space. 2) Do not consider detection of forms inside result pages.
Garvano et.al.[23]	1) Designed meta-search tools useful in searching over non-crawlable contents that are “hidden” behind search interfaces. 2) Content-summary construction technique	1) Query chosen only by using hierarchical categories as in Yahoo! and does not consider flat classification.
Bergholz et.al.[24]	1) Automatic identification of domain-specific Hidden Web resources. 2) Domain- specific crawling technique	1) Only deal with full text search forms. 2) Initialized with pre-classified documents and relevant keywords
Our Proposed Crawler	1) Use of Breadth-first search leads to wider coverage of web and achieves higher harvest rate. 2) By ranking collected sites and by focusing the crawling on a topic, it achieves more accurate results.	1) Crawling large amount of data causes time consumption.

Table 2: Comparison of Different Crawlers



ISSN(Online) : 2320-9801  
ISSN (Print) : 2320-9798

# International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Vol. 4, Issue 4, April 2016

At the end of this comparison we can conclude that the despite of having some limitations, proposed crawler is feasible to implement as well as more efficient than others. We still need to find a way to reduce the time consumption.

## V. CONCLUSION AND FUTURE WORK

This Paper surveys on different methods of crawling. Previous systems face many challenges such as efficiency, end-to-end delay, quality of link, failure to find the deep websites as they are unregistered with any crawler, scattered and dynamic. Thus, we proposed an effective and adaptive framework for retrieving deep web pages, namely Smart web Crawler. This approach accomplishes wide spread coverage of deep web and implements proficient crawling technique. We have used a focused crawler that comprises of two layers: Site-discovery and In-depth Crawling. It performs Site-discovery using Reverse Searching technique to fetch center pages from the known deep web pages and hence relevantly finds many data sources from several domains. Smart Web Crawler achieves more accuracy by correctly prioritizing the gathered sites and concentrating on the given domain. The In-depth crawling layer uses Smart Learning to perform search within a site and design a link hierarchy for avoiding biasness towards certain directories of a website for wider coverage of web. The comparative study shows the effectiveness of the proposed crawler which achieves higher harvest rates and wider coverage than other crawlers combining the Pre-querying and Post-querying approaches.

In future, we can incorporate this crawler with machine learning techniques to act and think like humans. This will enable the crawler to give results based upon the context. In addition to this, there is a need to expand the deep web dataset to crawl more number of websites giving maximum possible search results.

## ACKNOWLEDGEMENT

We would like to take this opportunity to express sincere gratitude and deep regards to our mentor and guide Prof. Ramachandra V. Pujeri, for his valuable feedback and constant encouragement throughout the process of writing this paper. Working under him was an extremely knowledgeable experience for us. We would like to thank our college, MIT College of Engineering (Pune), for providing us the necessary infrastructure and facilities. We would also like to give our gratitude to all the friends and teachers who helped us in every possible way and without whom this survey would be incomplete.

## REFERENCES

- [1] Anjum Asma and Gihan Nagib, 'Energy Efficient Routing Algorithms for Mobile Ad Hoc Networks—A Survey', International Journal of Emerging Trends & Technology in Computer Science, Vol.3, Issue 1, pp. 218-223, Rabia Iffat, Lalitha K. Sami, "Understanding the Deep Web", www.webpages.uidaho.edu/iffatsami.html
- [2] Dhiraj Khurana, Satish Kumar, "Web Crawler: A Review", International Journal of Computer Science & Management Studies (IJCSMS), Vol. 12, Issue 01, January 2012, ISSN (Online): 2231–5268.
- [3] Mini Singh Ahuja, Dr. Jatinder Singh Bal, Varnica, "Web Crawler: Extracting the Web Data", International Journal of Computer Trends and Technology (IJCTT) – volume 13 <http://www.ijcttjournal.org>
- [4] Devendra Hapase, Prof. M. D. Ingle, "Survey on Crawler for Deep-Web Interfaces", International Journal of Science and Research (IJSR) ISSN (Online): 2319-7064.
- [5] Olston and M. Najork, "Web Crawling, Foundations and Trends in Information Retrieval", vol. 4, No. 3, pp. 175–246, 2010.
- [6] Soumen Chakrabarti, Martin van den Berg 2, Byron Domc, "Focused crawling: a new approach to topic-specific Web resource discovery", Published by Elsevier Science B.V. All rights reserved in 1999.
- [7] Kevin Chen-Chuan Chang, Bin He, Chengkai Li, Mitesh Patel, and Zhen Zhang, "Structured databases on the web: Observations and implications", ACM SIGMOD Record, 33(3):61–70, 2004.
- [8] Sriram Raghavan and Hector Garcia-Molina, "Crawling the hidden web", In Proceedings of the 27th International Conference on Very Large Data Bases, pages 129–138, 2000.
- [9] Jayant Madhavan, Shawn R. Jeffery, Shirley Cohen, Xin Dong, David Ko, Cong Yu, and Alon Halevy, "Web-scale data integration: You can only afford to pay as you go", In Proceedings of CIDR, pages 342–350, 2007.
- [10] M. Franklin, A. Halevy, and D. Maier, "From databases to dataspace: A new abstraction for information management", *Sigmod Record*, 34(4):27–33, 2005.
- [11] Jared Cope, Nick Craswell, and David Hawking, "Automated discovery of search interfaces on the web", In Proceedings of the 14th Australasian database conference-Volume 17, pages 181–189. Australian Computer Society, Inc., 2003.
- [12] Thomas Kabisch, Eduard C. Dragut, Clement Yu, and Ulf Leser, "Deep web integration with VisQi", Proceedings of the VLDB Endowment, 3(1-2):1613–1616, 2010.



ISSN(Online) : 2320-9801  
ISSN (Print) : 2320-9798

# International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Vol. 4, Issue 4, April 2016

- [13] Luciano Barbosa, Juliana Freire, "Searching for Hidden Web Databases", Eighth International Workshop on the Web and Databases (WebDB 2005), June 16-17, 2005.
- [14] Kevin Chen-Chuan Chang, Bin He, and Zhen Zhang, "Toward large scale integration: Building a metaquerier over databases on the web", In *CIDR*, pages 44–55, 2005.
- [15] Denis Shestakov, "Databases on the web: national web domain survey", In Proceedings of the 15th Symposium on International Database Engineering & Applications, pages 179–184. ACM, 2011.
- [16] Denis Shestakov and Tapio Salakoski, "Host-ip clustering technique for deep web characterization", In Proceedings of the 12th International Asia-Pacific Web Conference (APWEB), pages 378–380. IEEE, 2010.
- [17] Denis Shestakov and Tapio Salakoski, "On estimating the scale of national deep web", In Database and Expert Systems Applications, pages 780–789. Springer, 2007.
- [18] Shestakov Denis, "On building a search interface discovery system", In Proceedings of the 2nd international conference on Resource discovery, pages 81–93, Lyon France, 2010. Springer.
- [19] Vinay Kancharla, "A Smart Web Crawler for a Concept Based Semantic Search Engine", San José State University, Master's Theses and Graduate Research, 2014.
- [20] Feng Zhao, Jingyu Zhou, Chang Nie, Heqing Huang, Hai Jin, "Smart Crawler: A Two-stage Crawler for Efficiently Harvesting Deep-Web Interfaces", IEEE Transactions on Services Computing Volume: PP Year: 2015.
- [21] Sriram Raghavan, Hector Garcia-Molina, "Crawling the Hidden Web", Computer Science Department Stanford University Stanford, CA 94305, USA.
- [22] Stephen W. Liddle, David W. Embley, Del T. Scott, and Sai Ho Yau, "Extracting Data behind Web Forms", Brigham Young University, Provo, UT 84602, USA.
- [23] Panagiotis G. Ipeirotis Luis Gravano, "Distributed Search over the Hidden Web: Hierarchical Database Sampling and Selection", Columbia University.
- [24] André Bergholz, Boris Chidlovskii, "Crawling for Domain-Specific Hidden Web Resources", Xerox Research Centre Europe 6 chemin de Maupertuis, 38240 Meylan, France.

## BIOGRAPHY

**Nimisha Jain** is a student in the Computer Engineering Department, MIT College of Engineering, Savitribai Phule Pune University. She is currently pursuing her Bachelor of Engineering.

**Pragya Sharma** is a student in the Computer Engineering Department, MIT College of Engineering, Savitribai Phule Pune University. She is currently pursuing her Bachelor of Engineering.

**Saloni Poddar** is a student in the Computer Engineering Department, MIT College of Engineering, Savitribai Phule Pune University. She is currently pursuing her Bachelor of Engineering.

**Shikha Rani** is a student in the Computer Engineering Department, MIT College of Engineering, Savitribai Phule Pune University. She is currently pursuing her Bachelor of Engineering.