# An Intelligent Crawler

[1] Mrugnayani Sharma, [2] Padmapani P. Tribhuvan
[1] PG Student, [2] Assistant Professor Department of Computer Science and Engineering
[1][2] Deogiri Institute of Engineering and Management,Studies,(DIEMS), Aurangabad, India

*Abstract:* A web crawler is a software program or programmed script that browses the world extensive web in a systematic, automated manner. Web crawler peregrinates from web page to page via the making use of the graphical structure of the internet pages. Such programs are additionally kenned as robots, spiders, and worms. In this system explained further, Data mining algorithms were used to introduce intelligence into the crawler. A statistical analysis of the performance of intelligent crawler is presented in this work. While introducing crawler intelligence, data mining algorithm plays an important role. The main objective is to develop an intelligent crawler to serve the purpose of web-indexing which helps in gathering relevant information from over the Internet with the help of search engines. The proposed intelligent crawler must perform crawling in minimum time with a maximum number of results.

*Keywords:* - Crawler, web indexing, statistical analysis.

## I. INTRODUCTION

Over the past decade, the web has grown exponentially, resulting in the prelude of the massive amount of data in the virtual world at every instant. Consequently, the conventional crawling strategy is eventually becoming inefficient in collecting and indexing web data. Thus intelligent crawlers must be developed and used to outperform the ever increasing Internet. The crawler is a multi-threaded bot that runs concurrently to serve the purpose of web-indexing which helps in gathering relevant information from over the Internet. This index is utilized by search engines, digital libraries, p2p communication, competitive perspicacity and many other industries. We are interested in introducing intelligent crawler which will perform crawling efficiently. Here the crawler is selective about the pages fetched and the links it will follow. This selectivity is based on the interest of the topic of the user thus at each step the crawler has to make a decision whether the next link will help to gather the content of interest. Other factors like a particular topic, the information it had already gathered also affect the efficiency and performance of the crawler [1]. While introducing perspicacity, two major approaches dominate the decisions made by the crawler. The first approach decides its crawling strategy by probing for the next best link amongst all links it can peregrinate whereas the second approach computes the benefit of peregrinating to all links and ranks them, which is utilized to decide the next link. The main objective is to develop perspicacity in crawler to accommodate the purport of web-indexing which avails in amassing pertinent information from over the Internet with the avail of search engines. The proposed perspicacious crawler must perform crawling in minimum time with a maximum number of results. As web crawler browses the World Wide Web in a methodical, automated manner, an astute web crawling strategy is to be implemented for boosting up the performance of crawling. The keenly intellective crawler must perform crawling in minimum time with maximum number of URLs crawled as a result [1, 2].

## II. LITERATURE SURVEY

TYPES OF CRAWLERS
### A. Parallel Crawlers
As the web grows in size, it becomes quite difficult or almost impossible to crawl the whole web by a single instance of a crawler. Therefore multiple processes are executed in parallel by search engines to cover the whole WWW. This type of crawler is referred to as a parallel crawler [3]. It consists of multiple crawling processes each of which performs the basic task of a single process crawler. The web pages are downloaded from the web and are stored locally. Afterwards, the URLs are extracted and their links are then followed.

### B. Focused Crawlers / Topical crawlers/ Topic driven crawlers
A focused crawler [3] has three main components a classifier that takes decisions on the relevancy of a page, a distiller decides the visit priorities and a crawler which downloads Webpages and is instructed by classifier and distiller module.

### C. Incremental Web Crawler
The incremental crawler [4,5] continuously crawls the web, revisiting pages periodically. During its continuous crawl, it may also purge some pages in the local collection, in order to make space for newly crawled pages. The crawler has following two goals:
• To keep the local collection fresh
• To improve quality of the local collection

### D. Hidden Web Crawler

Web crawlers generally crawl the web's dense tree structure called the publicly index able Web, i.e., the set of web pages reachable purely by following hypertext links. The surface web crawlers ignore search forms and pages that require authorization or prior registration. In particular, they ignore the huge amount of high-quality content "hidden" behind the search for. The Hidden web crawler [3, 4], called HiWE runs in a sequence of steps.

### III DEEP WEB CRAWLER'S FRAMEWORK [6]

The fundamental activities of a profound web crawler are like those of other conventional crawlers [6, 7].In Figure 1 the flowchart demonstrates the normal crawler circle, comprising of URL choice, page recovery, and page preparing to extricate joins. Note that customary crawlers don't recognize pages with and without shapes. The deep web crawler's execution sequence contains additional steps for pages on which forms are detected. Specifically, deep web crawler performs the following sequence of actions for each form on a page:

Step 1 Parse and process the form to build an internal representation, based on the model outlined in Section2. (Form Analysis)

Step 2 Generate the best (untried) value assignment and submit a completed form using that assignment.(Value assignment and submission)

Step 3 Analyze the response page to check if the submission yielded valid search results or if there were no matches. This feedback could be used to tune the value assignments in step 2.( Response Analysis)

Step 4 If the response page contains hypertext links, these are followed immediately (except for links that have already been visited or added to the queue) and recursively, to some pre-specified depth. Note that we could as well have added the links in the response page to the URL queue. However, for ease of implementation, in deep web crawler, we chose tonavigate the response pages immediately and that too, only up to a depth of 1.(Response Navigation)Steps 2, 3, and 4 are executed repeatedly, using different value assignments during each iteration. The sequence of value assignments is generated using the model.
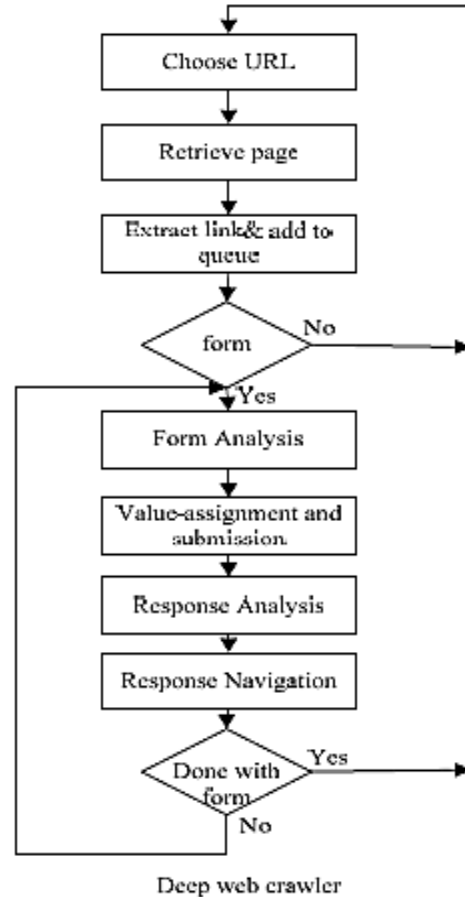


*Figure 1: Deep Web Crawler Loop*

The flowchart in fig. 1 illustrates the complete architecture of the deep web crawler. It includes six basic functional modules and two internal crawler data structures. The basic crawler data structure is the URL List. It contains all the URLs that the crawler has discovered so far. When starting up the crawler, the URL List is initialized to a seed set of URLs.

### IV.FUNCTIONS OF A WEB CRAWLER

The web searching process has two main components: offline and online [8]. The offline part is periodically performed by the search engine and it is used for building a collection of pages that will be later converted into a search index. The online part is executed each time an interrogation is performed by user. It uses the index for selecting documents that will later be sorted depending on estimation on their relevance with regard to the user's requirements. A schematic representation of this process is shown in figure 2.As web pages have different formats, the first stage for indexing web pages is represented by the extraction of a set of keywords.
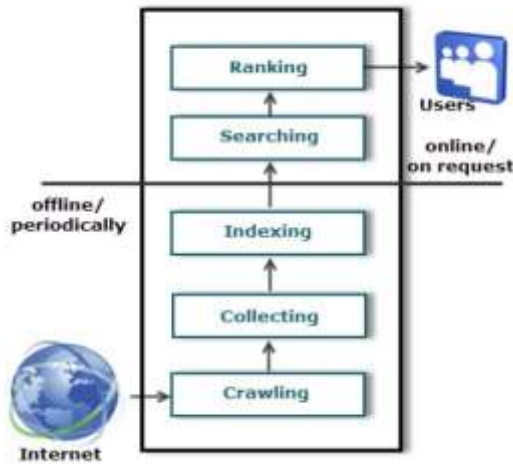
*Figure 2: General structure of a web searching process*

## V. ARCHITECTURE OF SMART CRAWLER FOR DEEP WEB INTERFACES

Smart Crawler consists of two main stages First is Site Locating and Second is In-site exploring [10]. The figure below shows the architecture of the proposed system.

Stage 1: Site locating– In Site locating stage the smart crawler performs the operation to find out the relevant sites related to the fired query. It has a number of steps involved to give the final result of this stage.

1) Seed Sites: It is the initial stage of the architecture. Here, seed sites are the candidate sites which are given to start crawling. It begins with the following URL of the query and explores other pages and other domains.
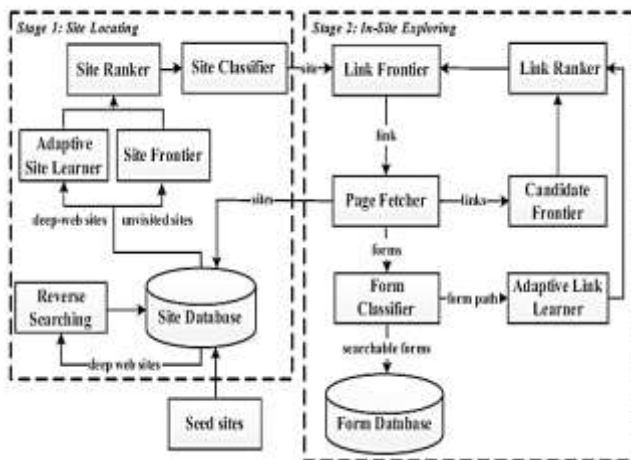


*Figure 3: Architecture of smart crawler for deep web interfaces*

2) **Reverse searching**: Pages with high rank and links to many other pages is called as a center page of the site. Some threshold is defined for seed sites, if a number of visited sited is less than the threshold then Reverse Searching is performed to know the center pages of the known deep web sites. Feed these pages back to the site database. The randomly picked site uses general search engine facility to find center pages and other relevant sites. Smart crawler first extract links on the page then download these pages and analyze these pages to decide whether the links are relevant or not. Following algorithm is used for reverse searching:

**Algorithm**
Input: seed sites and harvested deep websites.
Output: relevant sites.
1 while # of candidate sites less than a threshold do
2 // pick a deep website
3 site = getDeepWebSite(siteDatabase,seedSites)
4 resultP age = reverseSearch(site)
5 links = extractLinks(resultP age)
6 for each link in links do
7 page = downloadPage(link)
8 relevant = classify(page)
9 if relevant then
10 relevantSites = extractUnvisitedSite(page)
11 Output relevantSites
12 end if
13 endfor-each
14 end while

3) Incremental site Prioritizing: Incremental site prioritizing is used to achieve broad coverage on websites. It records the learned pattern of deep sites and forms the path for crawling. Basic knowledge is used to initialize both rankers such as site ranked and link ranker. Unvisited sites given to site frontier later prioritize by site ranked and added to the list fetched site. Two queues are used to classify out of site links such high priority queue and low priority queue respectively. High priority queue consist of out of site links which are classifieds relevant and judge by form classifier and low priority queue consist of links that are only judged as relevant. Algorithm for Incremental site Prioritizing is given below:

**Algorithm:**
Input: Site Frontier.
Output: searchable forms and out-of-site links.
1 HQueue=SiteFrontier.CreateQueue(HighPriority)
2 LQueue=SiteFrontier.CreateQueue(LowPriority)
3 while siteFrontier is not empty do
4 if HQueue is empty then
5 HQueue.addAll(LQueue)
6 LQueue.clear()
7 end

8 site = HQueue.poll()
9 relevant = classifySite(site)
10 if relevant then
11 performInSiteExploring(site)
12 Output forms and OutOfSiteLinks
13 siteRanker.rank(OutOfSiteLinks)
14 if forms is not empty then
15 HQueue.add (OutOfSiteLinks)
16 end
17 else
18 LQueue.add(OutOfSiteLinks)
19 end
20 end
21 end

**3) Site Frontier**: Site Frontier fetches the homepage URLsfrom the site database which is further ranked by Site Ranker to prioritize the highly relevant sites. Finding out-of-site links from visited web pages may not be enough for the Site Frontier.

**4) Adaptive link learner**: Site ranker and link ranker are controlled by Adaptive link learner. The feature space is decided for deep websites and links known as FSS and FSLrespectively. The Site Ranker is improved during crawling by an Adaptive Site Learner, which adaptively learns from features of deep-web sites (websites containing one or more searchable forms) found. The Link Ranker is adaptively improved by an Adaptive Link Learner, which learns from theURL path leading to relevant forms.

5) Site Ranker: Site ranker is used to rank unvisited site from the deep website. There are two parameters that are used for ranking mechanism are Site Similarity and Site Frequency. Site Similarity depends on the topic similarity between the known deep site and new site. Site Frequency is the occurrence of the site in another website.

6) Site Classifier: The high priority queue is for out-of-site links that are classified as relevant by Site Classifier and are judged by Form Classifier to contain searchable forms. If the site is the judge as atopic relevant then site crawling process is started otherwise the new site is picked from site frontier.

**B. Stage 2:** In-Site Exploring –After finding most relevant sites in stage 1 stage 2 perform the in-site exploration to find searchable forms.

**1) Link Frontier**: Link frontier takes sites as inputs which are classified by site classifier. Link frontier mainly works for finding links withincenter pages. Criteria for stopping early are given as Crawling Strategies: Mainly two crawling strategies are present Stop early and Balance link prioritizing.

Stop Early:
SC1: when reached maximum depth.
SC2: maximum crawling pages in each depth are reached.
SC3: Predefined numbers of forms are found at each depth.
SC4: No searchable forms till threshold value.
Balance link prioritizing: Here, link tree is constructed. The rootnode is the selected site and internal leaf node is each directory present on the website.

**2) Link Ranker**: Link Ranker prioritizes links so that SmartCrawler can quickly discover searchable forms. A high relevance score is given to a link that is most similar to links that directly point to pages with searchable forms.

3) Page Fetcher: Page Fetcher directly fetches out a center page of the website.

**4) Candidate Frontier:** The links in web pages are extracted into Candidate Frontier. The working of candidate frontier is similar as site frontier.

**5) Form Classifier:** Form classifier filters out non-searchable and irrelevant forms. The HIFI strategy is used to filter forms. HIFI consists of two classifiers, Searchable form classifier (SFC)and domain-specific form classifier(DSFC). SFC is domain independent and it filters out the non-searchable forms. It uses C4.5 algorithm for classification. DSFC is domain dependent and finds out the domain dependent form. Discuses Support vector machine.

6) Adaptive Link Learner: The Link Ranker is adaptively enhanced by an Adaptive Link Learner, which gains from the URL way prompting applicable structures.

**7) Form Database**: Form database contains a collection of sites; it collects all data which got input from Form Classifier.
At long last the outcome got is the most significant structures are acquired in profound web interfaces which are the coveted aftereffect of the proposed framework.

## VI. CONCLUSION

The savvy crawler is a compelling structure for the profound web. In this approach, we have accomplished wide scope for profound web and productive creeping. Savvy crawler s two-phase crawler comprising site situating by the turn around seeking with focus pages and on-site investigating comprises versatile connection positioning and connection tree for the more extensive scope.

## VII. ACKNOWLEDGEMENTS

## REFERENCES

[1] https://en.wikipedia.org/wiki/Web_crawler

[2] AbhirajDarshakar, Crawler intelligence with Machine Learning and Data Mining integration, Pune Institute of Computer Technology, Katraj, Pune, India (ICCCA2015) ISBN:978-1-4799-8890-7/15/$31.00 ©2015 IEEE 849

[3] Shruti Sharma and Parul Gupta, The Anatomy of Web Crawlers ISBN:978-1-4799-8890-7/15/$31.00 ©2015 IEEE

[4] Cho, J. and Garcia-Molina, H. 2003. Estimating frequency of change.ACM Transactions on Internet Technology 3, 3 (August).

[5] Cho J and Hector Garcia-Molina, "The evolution of the Web and implications for an incremental crawler", Prc. Of VLDB Conf., 2000.

[6] Xiang Peisu, TianKe and Huang Qinzhen, A Framework of Deep Web Crawler.

[7] JUNGHOO C, HECTOR GM, and LAWRENCE P. Efficient crawling through URLordering. Proceedings of the Seventh

[8] MirelaPirnau, Considerations on the functions and importance of a web crawler, ECAI 2015 - International Conference – 7th Edition Electronics, Computers, and Artificial Intelligence 978-1-4673-6647-/15/$31.00©2015 IEEE

[9] Keerthi S. Shetty, SwarajBhat and Sanjay Singh, Symbolic Verification of Web Crawler Functionality and Its Properties, 2012 International Conference OnComputerCommunication and Informatics (ICCCI -2012), Jan.10–12,2012, Coimbatore, INDIA

[10] Feng Zhao, Jingyu Zhou, Chang Nie, Heqing Huang, Hai Jin, SmartCrawler: A Two-stage Crawler for Efficiently Harvesting Deep-Web Interfaces,DOI 10.1109/TSC.2015.2414931, IEEE Transactions on Services Computing.