

# Search Optimization Using Smart Crawler

Dr. Mohammed Abdul Waheed<sup>1</sup>, Ajayraj Reddy<sup>2</sup>

<sup>1</sup>Associate Professor, Department of Computer Science & Engineering,

<sup>2</sup>P.G.Student, Department of Computer Science & Engineering,

VTU PG Centre, Kalaburgi, Karnataka, India.

**Abstract**— As deep net grows at a really quick pace, there has been multiplied interest in techniques that facilitate efficiently find deep-web interfaces. However, because of the massive volume of net resources and also the dynamic nature of deep net, achieving wide coverage and high efficiency may be a difficult issue. We tend to propose a two-stage framework, specifically Advance Crawler (ACrawler), for efficient gathering deep net interfaces. Within the first stage, ACrawler performs site-based sorting out centre pages with the assistance of search engines, avoiding visiting an oversized variety of pages. To realize additional correct results for a targeted crawl, ACrawler ranks websites to order extremely relevant ones for a given topic. Within the second stage, ACrawler achieves quick in-site looking by excavating most relevant links with associate degree accommodative link-ranking.

**Keywords**— peer-to-peer, store-carry-forward, discover-predict-deliver, interaction, recommendation.

## I. INTRODUCTION

A Web Crawler (also known as a robot or a spider) is a system for the bulk downloading of web pages. Web crawlers are used for a variety of purposes. Most prominently, they are one of the main components of web search engines, systems that assemble a corpus of web pages, index them, and allow users to issue queries against the index and find the web pages that match the queries. A related use is web archiving (a service provided by e.g., the Internet archive [3]), where large sets of web pages are periodically collected and archived for posterity. A third use is web data mining, where web pages are analyzed for statistical properties, or where data analytics is performed on them (an example would be Attributor [5], a company that monitors the web for copyright and trademark infringements). Finally, web monitoring services allow their clients to submit standing queries, or triggers, and they continuously crawl the web and notify clients of pages that match those queries. The deep (or hidden) web refers to the contents lie behind searchable web interfaces that cannot be indexed by searching engines. Based on extrapolations from a study done at University of California, Berkeley, it is estimated that the deep web contains approximately 91,850 terabytes and the surface web is only about 167 terabytes in 2003 [1]. More recent studies estimated that 1.9 zettabytes were reached and 0.3 zettabytes were consumed worldwide in 2007 [2], [3]. An IDC report estimates that the total of all digital data created, replicated, and consumed will reach 6 zettabytes in 2014 [4].

A significant portion of this huge amount of data is estimated to be stored as structured or relational data in web databases — deep web makes 96% of all the content on the Internet, which is 500-550 times larger than the surface web [4], [3]. These data contain a vast amount of valuable information and entities such as Infomine [5], Clusty [3], Books In Print [4] may be interested in building an index of the deep web sources in a given domain (such as book). Because these entities cannot access the proprietary web indices of search engines (e.g., Google and Baidu).

## II. OBJECTIVES

- 1) The Objective is to record learned patterns of deep web sites and form paths for incremental crawling.
- 2) Ranks site URLs to prioritize potential deep sites of a given topic. To this end, two features, site similarity and site frequency, are considered for ranking.
- 3) Focused crawler consisting of two stages: efficient site locating and balanced in-site exploring. SmartCrawler performs site-based locating by reversely searching the known deep web sites for center pages, which can effectively find many data sources for sparse domains.
- 4) SmartCrawler has an adaptive learning strategy that updates and leverages information collected successfully during crawling.

## III. WEB CRAWLER

A web crawler (also known as a robot or a spider) is a system, a program that traverses the web for the purpose of bulk downloading of web pages in an automated manner. Web crawlers are prominently one of the main components of web search engines that assemble a corpus of web pages or creates a copy of all the visited pages, index them, and allow users to issue queries against the index, provide fast searches and find the web pages that match the queries. Interacting with hundreds of thousands of web servers and name servers, crawling is considered as the most fragile application since it is beyond the control of the system. Crawler follows very simple steps yet very effective work in maintenance, checking of the downloaded links and also the validation of HTML

codes as follows It starts with the list of URL\_s to visit, called seeds and downloads the web page.

#### IV. PROBLEM DEFINITION

An effective deep web harvesting framework, namely SmartCrawler, is for achieving both wide coverage and high efficiency for a focused crawler. Based on the observation that deep websites usually contain a few searchable forms and most of them are within a depth of three. Our crawler is divided into two stages: site locating and in-site exploring. The site locating stage helps achieve wide coverage of sites for a focused crawler, and the in-site exploring stage can efficiently perform searches for web forms within a site.

#### V. RELATED WORK

Existing strategies were dealing with creation of a single profile per user, but conflict occurs when user's interest varies for the same query Eg. When a user is interested in banking exams in query —bankl may be slightly interested in accounts of money bank where not at all interested in blood bank. At such time conflict occurs so we are dealing with negative preferences to obtain the fine grain between the interested results and not interested. Consider following two aspects:

##### 1) Document-Based method:

These methods aim at capturing users's clicking and browsing behavior. It deals with click through data from the user i.e. the documents user has clicked on. Click through data in search engines can be thought of as triplets (q, r, c)

Where,

q = query

r = ranking

c = set of links clicked by user.

##### 2) Concept-based methods:

These methods aim at capturing user's conceptual needs. User's browsed documents and search histories. User profiles are used to represent user's interests and to infer their intentions for new queries.

#### Disadvantages of Existing System:

- 1) Deep-web interfaces.
- 2) Achieving wide coverage and high efficiency is a challenging issue.

#### VI. PROPOSED SYSTEM MECHANISM

To efficiently and effectively discover deep web data sources, SmartCrawler is designed with two stage architecture, site locating and in-site exploring, as shown in Figure 1. The first site locating stage finds the most relevant site for a given topic, and then the second in-site exploring stage uncovers searchable forms from the site. Specifically, the site locating

stage starts with a seedset of sites in a site database. Seeds sites are candidate sites given for SmartCrawler to start crawling, which begins by following URLs from chosen seed sites to explore other pages and other domains. When the number of unvisited URLs in the database is less than a threshold during the crawling process, SmartCrawler performs [reverse searching] of known deep websites for center pages (highly ranked pages that have many links to other domains) and feeds these pages back to the site database. Site Frontier fetches homepage URLs from the site database, we going to rank the relevant information.

#### VII. SYSTEM ARCHITECTURE: (TWO STAGE ARCHITECTURE)

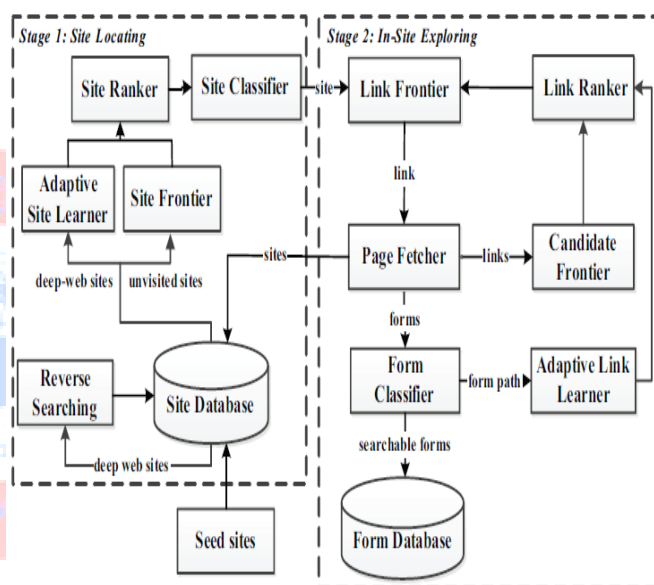


Fig.1 System Architecture: (Two stage Architecture)

To efficiently and effectively discover deep web data sources, SmartCrawler is designed with two stage architecture, site locating and in-site exploring, as shown in Figure. The first site locating stage finds the most relevant site for a given topic, and then the second in-site exploring stage uncovers searchable forms from the site. Specifically, the site locating stage starts with a seed set of sites in a site database. Seeds sites are candidate sites given for SmartCrawler to start crawling, which begins by following URLs from chosen seed sites to explore other pages and other domains. When the number of unvisited URLs in the database is less than a threshold during the crawling process, SmartCrawler performs [reverse searching] of known deep websites for center pages (highly ranked pages that have many links to other domains) and feeds these pages back to the site database.

#### IV. IMPLEMENTATION

##### A. ALGORITHMS & TECHNIQUES USED

###### Algorithm 1: Reverse searching for more sites.

**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 foreach 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
13 end
14 end
```

###### Algorithm 2: Incremental Site Prioritizing.

**Input :** siteFrontier

**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
```

##### B. MODULE INFORMATION

###### 1. Two-stage crawler

It is challenging to locate the deep web databases, because they are not registered with any search engines, are usually sparsely distributed, and keep constantly changing. To address this problem, previous work has proposed two types of crawlers, generic crawlers and focused crawlers. Generic crawlers fetch all searchable forms and cannot focus on a

specific topic. Focused crawlers such as Form-Focused Crawler (FFC) and Adaptive Crawler for Hidden-web Entries (ACHE) can automatically search online databases on a specific topic. FFC is designed with link, page, and form classifiers for focused crawling of web forms, and is extended by ACHE with additional components for form filtering and adaptive link learner. The link classifiers in these crawlers play a pivotal role in achieving higher crawling efficiency than the best-first crawler. However, these link classifiers are used to predict the distance to the page containing searchable forms, which is difficult to estimate, especially for the delayed benefit links (links eventually lead to pages with forms). As a result, the crawler can be inefficiently led to pages without targeted forms.

###### 2. Site Ranker

When combined with above stop-early policy. We solve this problem by prioritizing highly relevant links with link ranking. However, link ranking may introduce bias for highly relevant links in certain directories. Our solution is to build a link tree for a balanced link prioritizing. Figure 2 illustrates an example of a link tree constructed from the homepage of <http://www.abebooks.com>. Internal nodes of the tree represent directory paths. In this example, servlet directory is for dynamic request; books directory is for displaying different catalogs of books; Amdocs directory is for showing help information. Generally each directory usually represents one type of files on web servers and it is advantageous to visit links in different directories. For links that only differ in the query string part, we consider them as the same URL. Because links are often distributed unevenly in server directories, prioritizing links by the relevance can potentially bias toward some directories. For instance, the links under books might be assigned a high priority, because —bookl is an important feature word in the URL. Together with the fact that most links appear in the books directory, it is quite possible that links in other directories will not be chosen due to low relevance score. As a result, the crawler may miss searchable forms in those directories.

###### 3. Adaptive learning

Adaptive learning algorithm that performs online feature selection and uses these features to automatically construct link rankers. In the site locating stage, high relevant sites are prioritized and the crawling is focused on atopic using the contents of the root page of sites, achieving more accurate results. During the in site exploring stage, relevant links are prioritized for fast in-site searching. We have performed an extensive performance evaluation of Smart Crawler over real web data in 1representativedomains and compared with ACHE and site-based crawler. Our evaluation shows that our crawling framework is very effective, achieving substantially higher harvest rates than the state-of-the-art ACHE crawler. The results also show the effectiveness of the reverse searching and adaptive learning.

## V. CONCLUSIONS

We've shown that our approach achieves each wide coverage for deep net interfaces and maintains extremely efficient locomotion. Crawdy may be a targeted crawler consisting of 2 stages: efficient web site locating and balanced in-site exploring. Crawdy performs site-based locating by reversely looking the glorious deep websites for center pages, which may effectively find several knowledge sources for distributed domains. By ranking collected sites and by focusing the locomotion on a subject, Crawdy achieves additional correct results. The in-site exploring stage uses adaptive link-ranking to go looking at intervals a site; and that we style a link tree for eliminating bias toward sure directories of an internet site for wider coverage of web directories. Our experimental results on a representative set of domains show the effectiveness of the projected two-stage crawler that achieves higher harvest rates than different crawlers. In future work, we tend to arrange to mix pre-query and postquery approaches for classifying deep-web forms to more improve the accuracy of the shape classifier.

## REFERENCES

- [1] Peter Lyman and Hal R. Varian. How much information? 2003. Technical report, UC Berkeley, 2003.
- [2] Roger E. Bohn and James E. Short. How much information? 2009 report on american consumers. Technical report, University of California, San Diego, 2009.
- [3] Martin Hilbert. How much information is there in the "information society"? Significance, 9(4):8–12, 2012.
- [4] Idc worldwide predictions 2014: Battles for dominance – and survival–on the 3rd platform. <http://www.idc.com/research/Predictions14/index.jsp>, 2014.
- [5] Michael K. Bergman. White paper: The deep web: Surfacing hidden value. Journal of electronic publishing, 7(1), 2001.
- [6] Yeye He, Dong Xin, Venkatesh Ganti, Sriram Rajaraman, and Nirav Shah. Crawling deep web entity pages. In Proceedings of the sixth ACM international conference on Web search and data mining, pages 355–364. ACM, 2013.
- [7] Infomine. UC Riverside library. <http://lib-www.ucr.edu/>, 2014.
- [8] Clusty's searchable database dirctory. <http://www.clusty.com/>, 2009.
- [9] Booksinprint. Books in print and global books in print access. <http://booksinprint.com/>, 2015.
- [10] 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.
- [11] 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.
- [12] 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.
- [13] Denis Shestakov and Tapio Salakoski. On estimating the scale of national deep web. In Database and Expert Systems Applications, pages 780–789. Springer, 2007.
- [14] 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.
- [15] Luciano Barbosa and Juliana Freire. Searching for hidden-web databases. In WebDB, pages 1–6, 2005.
- [16] Luciano Barbosa and Juliana Freire. An adaptive crawler for locating hidden-web entry points. In Proceedings of the 16th international conference on World Wide Web, pages 441–450. ACM, 2007.
- [17] Soumen Chakrabarti, Martin Van den Berg, and Byron Dom. Focused crawling: a new approach to topic-specific web resource discovery. Computer Networks, 31(11):1623–1640, 1999.
- [18] Jayant Madhavan, David Ko, Lucja Kot, Vignesh Ganapathy, Alex Rasmussen, and Alon Halevy. Google's deep web crawl. Proceedings of the VLDB Endowment, 1(2):1241–1252, 2008.
- [19] Olston Christopher and Najork Marc. Web crawling. Foundations and Trends in Information Retrieval, 4(3):175–246, 2010.
- [20] Balakrishnan Raju and Kambhampati Subbarao. Sourcerank: 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, 2011.
- [21] Balakrishnan Raju, Kambhampati Subbarao, and Jha Manishkumar. 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, 2013.
- [22] Mustafa Emmre Dincturk, Guy vincent Jourdan, Gregor V. Bochmann, and Iosif Viorel Onut. A model-based approach for crawling rich internet applications. ACM Transactions on the Web, 8(3):Article 19, 1–39, 2014.
- [23] 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.
- [24] Wensheng Wu, Clement Yu, AnHai Doan, and Weiyi Meng. An interactive clustering-based approach to integrating source query interfaces on the deep web. In Proceedings of the 2004 ACM SIGMOD international conference on Management of data, pages 95–106. ACM, 2004.
- [25] Eduard C. Dragut, Thomas Kabisch, Clement Yu, and Ulf Leser. A hierarchical approach to model web query interfaces for web source integration. Proc. VLDB Endow., 2(1):325–336 August 2009.
- [26] 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.
- [27] Eduard C. Dragut, Weiyi Meng, and Clement Yu. Deep Web Query Interface Understanding and Integration. Synthesis Lectures on Data Management. Morgan & Claypool Publishers, 2012.
- [28] Andr'e Bergholz and Boris Childlovskii. Crawling for domainspecific hidden web resources. In Web Information Systems Engineering, 2003. WISE 2003. Proceedings of the Fourth International Conference on, pages 125–133. IEEE, 2003.
- [29] 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.
- [30] Cheng Sheng, Nan Zhang, Yufei Tao, and Xin Jin. Optimal algorithms for crawling a hidden database in the web. Proceedings of the VLDB Endowment, 5(11):1112–1123, 2012.
- [31] Panagiotis G Ipeirotis and Luis Gravano. Distributed search over the hidden web: Hierarchical database sampling and selection. In Proceedings of the 28th international conference on Very Large Data Bases, pages 394–405. VLDB Endowment, 2002.
- [32] Nilesh Dalvi, Ravi Kumar, Ashwin Machanavajjhala, and Vibhor Rastogi. Sampling hidden objects using nearest-neighbor oracles. In Proceedings of the 17th ACM SIGKDD international conference on Knowledge discovery and data mining, pages 1325–1333. ACM, 2011.
- [33] 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.



- [34] Brightplanet's searchable database directory. <http://www.completeplanet.com/>, 2001.
- [35] Mohamamdreza Khelghati, Djoerd Hiemstra, and Maurice Van Keulen. Deep web entity monitoring. In Proceedings of the 22nd international conference on World Wide Web companion, pages 377–382. International World Wide Web Conferences Steering Committee, 2013.
- [36] Soumen Chakrabarti, Kunal Punera, and Mallela Subramanyam. Accelerated focused crawling through online relevance feedback. In Proceedings of the 11th international conference on World Wide Web, pages 148–159, 2002.
- [37] Luciano Barbosa and Juliana Freire. Combining classifiers to identify online databases. In Proceedings of the 16th international conference on World Wide Web, pages 431–440. ACM, 2007.
- [38] 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.
- [39] Dumais Susan and Chen Hao. Hierarchical classification of Web content. In Proceedings of the 23rd Annual International ACM SIGIR conference on Research and Development in Information Retrieval, pages 256–263, Athens Greece, 2000.
- [40] Mark Hall, Eibe Frank, Geoffrey Holmes, Bernhard Pfahringer, Peter Reutemann, and Ian H. Witten. The weka data mining software: an update. SIGKDD Explorations Newsletter, 11(1):10–18, November 2009.
- [41] Open directory project. <http://www.dmoz.org/>, 2013.
- [42] The UIUC web integration repository. <http://metaquerier.cs.uiuc.edu/repository/>, 2003.

