

# Survey on Crawler for Harvesting Deep Web Interfaces

Pooja M. Taide<sup>1</sup> Vijaya Kamble<sup>2</sup>

<sup>1,2</sup>M.Tech Student

<sup>1,2</sup>Department of Computer Science & Engineering

<sup>1,2</sup>Gurunanak Institute of Technology, Nagpur, India

**Abstract**— Due to heavy usage of internet large amount of diverse data is spread over it which provides access to particular data or to search most relevant data. It is very challenging for search engine to fetch relevant data as per user's need and which consumes more time. So, to reduce large amount of time spend on searching most relevant data we proposed the "Smart Crawler". In this proposed crawler framework, which efficiently retrieves deep-web interfaces from large-scale sites, easily get the information which is stored in web databases. In the first stage, Smart Crawler performs site-based searching for center pages with the help of search engines, avoiding visiting a large number of pages. To achieve more accurate results for a focused crawl, Smart Crawler ranks websites to prioritize highly relevant ones for a given topic. In the second stage, Smart Crawler achieves fast in-site searching by excavating most relevant links with an adaptive link-ranking. To eliminate bias on visiting some highly relevant links in hidden web directories, we design a link tree data structure to achieve wider coverage for a website. Our experimental results on a set of representative domains show the agility and accuracy of our proposed crawler framework, which efficiently retrieves deep-web interfaces from large-scale sites and achieves higher harvest rates than other crawlers. We propose a two stages framework, namely Smart Crawler, for efficient harvesting deep web interfaces.

**Key words:** Deep Web, Crawler, Feature Selection, Ranking, Adaptive Learning

## I. INTRODUCTION

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 up about 96% of all the content on the Internet, which is 500-550 times larger than the surface web [5], [6].

These data contain a vast amount of valuable information and entities such as Info mine, Clusty, Books In Print 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), there is a need for an efficient crawler that is able to accurately and quickly explore the deep web databases.

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.

In this paper, we propose an effective deep web harvesting framework, namely Smart Crawler, 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. Our main contributions are:

We propose a novel two-stage framework to address the problem of searching for hidden-web resources. Our site locating technique employs a reverse searching technique (e.g., using Google's "link:" facility to get pages pointing to a given link) and incremental two-level site prioritizing technique for unearthing relevant sites, achieving more data sources. During the in-site exploring stage, we design a link tree for balanced link prioritizing, eliminating bias toward webpages in popular directories.

We propose an 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 a topic using the contents of the root page of sites, achieving more accurate results. During the insite exploring stage, relevant links are prioritized for fast in-site searching.

## II. LITERATURE SURVEY

For our topic smart crawler we have done a survey on some of the IEEE and other standard papers. These papers and their description are as follows:

A. 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.*

From this research paper we got information that, as deep net grows at a really quick pace, there has been multiplied interest in techniques that facilitate efficiently and 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 Crawdy, for efficient gathering deep net interfaces. Within

the first stage, Crawdy 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, Crawdy ranks websites to order extremely relevant ones for a given topic. Within the second stage, Crawdy achieves quick in-site looking by excavating most relevant links with associate degree accommodative link-ranking.

*B. 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, 132, 2013.*

This paper helps us to know about the one immediate challenge in searching the deep web databases is source selection i.e. selecting the most relevant web databases for answering a given query. The existing methods of database selection (both text and relational databases) uses relevance measures based on the similarity with the queries for the quality assessment of the sources. Existing methods have two deficiencies for applying to the open collections like the deep web. First is that the methods are agnostic to the correctness (trustworthiness) of the sources. Secondly, since the existing measures are fully dependent on the query similarity, they do not consider the popularity of the results for computing the probability of relevance. Since number of sources provide their own answer sets to the same query in the deep web, the agreements between these answer sets are likely to be helpful in assessing the relevance and trustworthiness of the sources. We start with this hypothesis and compute the agreement between the sources using entity matching methods. Agreement is modeled as a graph with sources at the vertices. On this agreement graph source quality score namely SourceRank is calculated as the stationary visit probability of a random walk. Our evaluations on the online deep web sources show that the relevances of the sources selected by SourceRank is improved by 20-50per over the existing methods; and that SourceRank of a source reduces linearly with corruption levels. Also we demonstrate that SourceRank can be combined with Google Base ranking to improve the precision by 22-60per and to select sources better trusted by the users.

*C. Mustafa Emmre Dincturk, Guy vincent Jourdan, Gregor V.IEEE Transactions on Services Computing Volume: PP Year: 201514Bochmann, and Iosif Viorel Onut. A model-based approach for crawling rich internet applications. ACM Transactions on the Web, 8(3):Article 19, 139, 2014.*

This paper gives knowledge about the New Web technologies, like AJAX, result in more responsive and interactive Web applications, sometimes called Rich Internet Applications (RIAs). Crawling techniques developed for traditional Web applications are not sufficient for crawling RIAs. The inability to crawl RIAs is a problem that needs to be addressed for at least making RIAs searchable and testable. We present a new methodology, called model-based crawling, that can be used as a basis to design efficient crawling strategies for RIAs. We illustrate model-based crawling with a sample strategy, called the hypercube strategy. The performances of our model-based crawling strategies are compared against existing standard crawling

strategies, including breadth-first, depth first, and a greedy strategy. Experimental results show that our model-based crawling approach is significantly more efficient than these standard strategies.

### III. RELATED WORK

To leverage the large volume information buried in deep web, previous work has proposed a number of techniques and tools, including deep web understanding and integration, hidden web crawlers [6], and deep web samplers. For all these approaches, the ability to crawl deep web is a key challenge. Olston and Najork systematically present that crawling deep web has three steps: locating deep web content sources, selecting relevant sources and extracting underlying content. Following their statement, we discuss the two steps closely related to our work as below.

Locating deep web content sources. A recent study shows that the harvest rate of deep web is low — only 647,000 distinct web forms were found by sampling 25 million pages from the Google index (about 2.5%). Generic crawlers are mainly developed for characterizing deep web and directory construction of deep web resources, that do not limit search on a specific topic, but attempt to fetch all searchable forms[5]. The Database Crawler in the MetaQuerier is designed for automatically discovering query interfaces. Database Crawler first finds root pages by an IP-based sampling, and then performs shallow crawling to crawl pages within a web server starting from a given root page. The IPbased sampling ignores the fact that one IP address may have several virtual hosts, thus missing many websites. To overcome the drawback of IPbased sampling in the Database Crawler, Denis et al. propose a stratified random sampling of hosts to characterize national deep web, using the Hostgraph provided by the Russian search engine Yandex. I-Crawler [5] combines pre-query and post-query approaches for classification of searchable forms.

Selecting relevant sources. Existing hidden web directories usually have low coverage for relevant online databases, which limits their ability in satisfying data access needs. Focused crawler is developed to visit links to pages of interest and avoid links to off-topic regions. Soumen et al. describe a best-first focused crawler, which uses a page classifier to guide the search.

The classifier learns to classify pages as topic-relevant or not and gives priority to links in topic relevant pages. However, a focused best-first crawler harvests only 94 movie search forms after crawling 100,000 movie related pages. An improvement to the best-first crawler is proposed in, where instead of following all links in relevant pages, the crawler used an additional classifier, the apprentice, to select the most promising links in a relevant page. The baseline classifier gives its choice as feedback so that the apprentice can learn the features of good links and prioritize links in the frontier.

Different from the crawling techniques and tools mentioned above, Smart Crawler is a domain-specific crawler for locating relevant deep web content sources. Smart Crawler targets at deep web interfaces and employs a two-stage design, which not only classifies sites in the first stage to filter out irrelevant websites, but also categorizes

searchable forms in the second stage. Instead of simply classifying links as relevant or not, Smart Crawler first ranks sites and then prioritizes links within a site with another ranker.

#### IV. PROPOSED SYSTEM

In this proposed system we use a Meta search engine; in Meta search engine the searching result is accumulated by using multiple search engines. Actually, it is good in finding the unique key word phrases, quotes, and Knowledge encompasses in the full text of web pages. And Search engines allow user to enter keywords and then examine this keyword in its table followed by database. We propose a novel two-stage framework to address the problem of searching for hidden-web Resources. Our site locating technique employs a reverse searching strategy(e.g. By using Google’s link: facility to get pages guiding to a given link) and incremental two-level site prioritizing technique to discover more relevant sites, achieving more data sources. During the in-site exploring stage, we have constructed a link tree for balanced link prioritizing, eliminating bias toward web pages in popular directories. We introduce an adaptive learning algorithm that performs online feature selection and uses these features to automatically create link rankers. In site locating stage, high relevant sites are arranged and the crawling is focused on a topic using the contents of the seed page of sites, achieving more accurate results. During the insite exploring stage, relevant links are arranged for fast in-site searching.

To efficiently and effectively discover deep web data sources, Web Crawler is designed with two stage architecture, site locating and in-site exploring, as shown in architecture diagram. 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 web crawler 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, Then "reverse searching " is perform by web crawler of known deep web sites 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 databases, which are ranked by Site Ranker to prioritize highly relevant sites. The Site Ranker is improved during crawling by an Adaptive Site Learner, which adaptively learns from features of deep-web sites (web sites containing one or more searchable forms) found. To achieve more accurate results for a focused crawl, Site Classifier categorizes URLs into relevant or irrelevant for a given topic according to the homepage content.

##### A. Two-Stage Architecture

To efficiently and effectively discover deep web data sources, Smart Crawler is designed with a two stage architecture, site locating and in-site exploring, as shown in Figure 1.

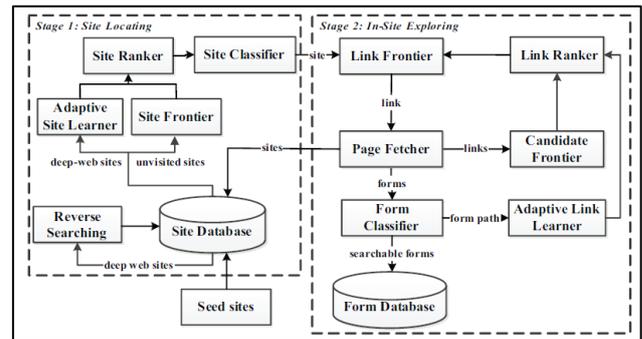


Fig. 1: The Two Stage Architecture of Smart Crawler

We propose a two-stage framework, namely Smart Crawler, for efficient harvesting deep web interfaces. In the first stage, Smart Crawler performs site-based searching for center pages with the help of search engines, avoiding visiting a large number of pages. To achieve more accurate results for a focused crawl, Smart Crawler ranks websites to prioritize highly relevant ones for a given topic. In the second stage, Smart Crawler achieves fast in-site searching by excavating most relevant links with an adaptive link-ranking. To eliminate bias on visiting some highly relevant links in hidden web directories, we design a link tree data structure to achieve wider coverage for a website. Our experimental results on a set of representative domains show the agility and accuracy of our proposed crawler framework, which efficiently retrieves deep-web interfaces from large-scale sites and achieves higher harvest rates than other crawlers. Propose an effective harvesting framework for deep-web interfaces, namely Smart-Crawler. We have shown that our approach achieves both wide coverage for deep web interfaces and maintains highly efficient crawling. Smart Crawler is a focused crawler consisting of two stages: efficient site locating and balanced in-site exploring. Smart Crawler performs site-based locating by reversely searching the known deep web sites for center pages, which can effectively and many data sources for sparse domains. By ranking collected sites and by focusing the crawling on a topic, Smart Crawler achieves more accurate results.

#### V. CONCLUSION

In this paper we have identified the different kind of general searching technique and Meta search engine strategy and by using this we have proposed an effective way of searching most relevant data for deep web interfaces.

#### 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] 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.
- [6] 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.
- [7] Idc worldwide predictions 2014: Battles for dominance and survival on the 3rd platform.
- [8] Michael K. Bergman. White paper: The deep web: Surfacing hidden value. Journal of electronic publishing, 7(1), 2001.
- [9] 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 datamining, pages 355–364. ACM, 2013.
- [10] Infomine. UC Riverside library. <http://lib-www.ucr.edu/>,
- [11] Clustys searchable database directory. <http://www.clusty.com/>, 2009.
- [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.

