# Smart Crawler for Efficient Deep-Web Harvesting

[1]Gayatri Dhumal, [2]Rucha Joshi, [3]Amarja Kaldate

[1]Gayatri Dhumal, Computer Engineering, PVPIT Bavdhan, Maharashtra, India
[2]Rucha Joshi, Computer Engineering, PVPIT Bavdhan, Maharashtra, India
[3]Amarja Kaldate, Computer Engineering, PVPIT Bavdhan, Maharashtra, India

*Abstract:* **As profound web developer at quick pace, there has been expanded enthusiasm for method that assists proficiently with finding profound web interfaces. Nonetheless, because of the extensive volume of web assets and the dynamic way of profound web, accomplishing wide scope and high productivity is a testing issue. We propose a two-stage structure, in particular Smart Crawler, for effective gathering profound web interfaces. In the first stage, Smart Crawler performs site-based hunting down focus pages with the assistance of web indexes, abstaining from going by countless. To accomplish more exact results for an engaged slither, Smart Crawler positions sites to organize profoundly pertinent ones for a given point. In the second stage, Smart Crawler accomplishes quick in-site excavating so as to see most significant connections with a versatile connection positioning. To dispense with inclination on going by some exceedingly significant connections in shrouded web indexes, we outline a connection tree information structure to accomplish more extensive scope for a site. Our test results on an arrangement of delegate areas demonstrate the readiness and precision of our proposed crawler structure, which effectively recovers profound web interfaces from huge scale destinations and accomplishes higher harvest rates than different crawlers.**

*Keywords:* **Deep web, two-stage crawler, feature selection, ranking, adaptive learning.**

## 1. INTRODUCTION

The profound (or shrouded) web alludes to the substance lie behind searchable web interfaces that can't be listed via looking motors. In light of extrapolations from a study done at University of California, Berkeley, it is evaluated that the profound web contains pretty nearly 91,850 terabytes and the surface web is just around 167 terabytes in 2003. Later studies evaluated that 1.9 zetta bytes were come to and 0.3 zetta bytes were expended worldwide in 2007. An IDC report assesses that the aggregate of all advanced information made, recreated, and expended will achieve 6 zetta bytes in 2014. A critical segment of this tremendous measure of information is evaluated to be put away as organized or social information in web databases - profound web makes up around 96% of all the substance on the Internet, which is 500-550 times bigger than the surface web. These information contain an inconceivable measure of important data and elements, for example, Info mine, Cluster, Books In Print may be keen on building a list of the profound web sources in a given area, (for example, book). Since these elements can't get to the restrictive web files of web crawlers, there is a requirement for an effective crawler that has the capacity precisely and rapidly investigates the profound web database.

It is trying to find the profound web databases, in light of the fact that they are not enlisted with any web indexes, are typically scantily conveyed, and keep continually evolving. To address this issue, past work has proposed two sorts of crawlers, nonexclusive crawlers and centered crawlers. Nonexclusive crawlers, get every single searchable structure and can't concentrate on a particular subject. Centered crawlers, for example, Form-Focused Crawler (FFC) and Adaptive Crawler for Hidden-web Entries (ACHE) can naturally seek online databases on a particular theme. FFC is outlined with connection, page, and structure classifiers for centered slithering of web structures, and is reached out by ACHE with

extra segments for structure separating and versatile connection learner. The connection classifiers in these crawlers assume a crucial part in accomplishing higher slithering proficiency than the best-first crawler. Notwithstanding, these connection classifiers are utilized to anticipate the separation to the page containing searchable structures, which is hard to assess, particularly for the deferred advantage connections (interfaces in the long run lead to pages with structures). Therefore, the crawler can be wastefully prompted pages without focused on structures.

In this paper, 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:

• This paper proposes 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 web pages 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 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.

## 2. LITERATURE SURVEY

**1) Host-ip clustering technique for deep web characterization:**

**AUTHORS:** Denis Shestakov and Tapio Salakoski.

A huge portion of today's web pages filled with information from myriads of online database. This part of the Web, known as the deep Web, is to date relatively unexplored and even major characteristics such as number of searchable databases on the Web is somewhat disputable. In this paper, we are aimed at more accurate estimation of main parameters of the deep Web by sampling one national web domain. We propose the Host-IP clustering sampling technique that addresses drawbacks of existing approaches to characterize the deep Web and report our findings based on the survey of Russian Web conducted in September 2006. Obtained estimates together with a proposed sampling method could be useful for further studies to handle data in the deep Web.

**2) Searching for hidden-web databases:**

**AUTHORS:** Luciano Barbosa and Juliana Freire.

Recently, there has been increased interest in the retrieval and integration of hidden-Web data with a view to leverage

High-quality information available in online database Although previous works have addressed many aspects of the actual integration, including matching form schemata and automatically filling out forms, the problem of locating relevant data sources has been largely overlooked.

Given the dynamic nature of the Web, where data sources are constantly changing, it is crucial to automatically discover the sere sources. However, considering the number of documents on the Web (Google already indexes over 8 billion documents), automatically finding tens, hundreds or even thousands of forms that are relevant to the integration task is really like looking for a few needles in a haystack. Besides, since the vocabulary and structure of forms for a given domain are unknown until the forms are actually found, it is hard to define exactly what to look for. We propose a new crawling strategy to automatically locate hidden-Web databases which aims to achieve a balance between the two conflicting requirements of this problem: the need to perform a broad search while at the same time avoiding the need to crawl a large number of irrelevant pages. The proposed strategy does that by focusing the crawl on a given topic; by judiciously choosing links to follow within a topic that are more likely to lead to pages that contain forms; and by employing appropriate stopping criteria. We describe the algorithms underlying this strategy and an experimental evaluation which shows that our approach is both effective and efficient, leading to larger numbers of forms retrieved as a function of the number of pages visited than other crawlers.

**ISSN 2350-1022**

**International Journal of Recent Research in Mathematics Computer Science and Information Technology**
Vol. 3, Issue 1, pp: (43-50) Month: April 2016 – September 2016, Available at: **www.paperpublications.org**

**3) Crawling for domain specific hidden web resources:**

**AUTHORS:** Andr´eBergholz and Boris Childlovskii.

The Hidden Web, the part of the Web that remains unavailable for standard crawlers, has become an in important tpic during recent years. Its size is estimated to 400 to 500 times larger than that of the publicly index able Web (PIW). Furthermore, the information on the hidden Web is assumed to be more structured, because it is usually stored in databases. In this paper, we describe a crawler which starting from the PIW finds entry points into the hidden Web. The crawler is domain-specific and is initialized with pre-classified documents and relevant keywords. We describe our approach to the automatic identification of Hidden Web resources among encountered HTML forms. We conduct a series of experiments using the top-level categories in the Google directory and report our analysis of the discovered Hidden Web resources.

**4) Crawling the hidden web:**

**AUTHORS:** Sriram Raghavan and Hector Garcia-Molina.

Current-day crawlers retrieve content only from the publicly index able Web, i.e., the set of Web pages reachable purely by following hypertext links, ignoring search forms and pages that require authorization or prior registration. In particular, they ignore the tremendous amount of high quality content ``hidden" behind search forms, in large searchable electronic databases. In this paper, we address the problem of designing a crawler capable of extracting content from this hidden Web. We introduce a generic operational model of a hidden Web crawler and describe how this model is realized in HIWE (Hidden Web Exposer), a prototype crawler built at Stanford. We introduce a new Layout-based Information Extraction Technique (LITE) and demonstrate its use in automatically extracting semantic information from search forms and response pages. We also present results from experiments conducted to test and validate our techniques.
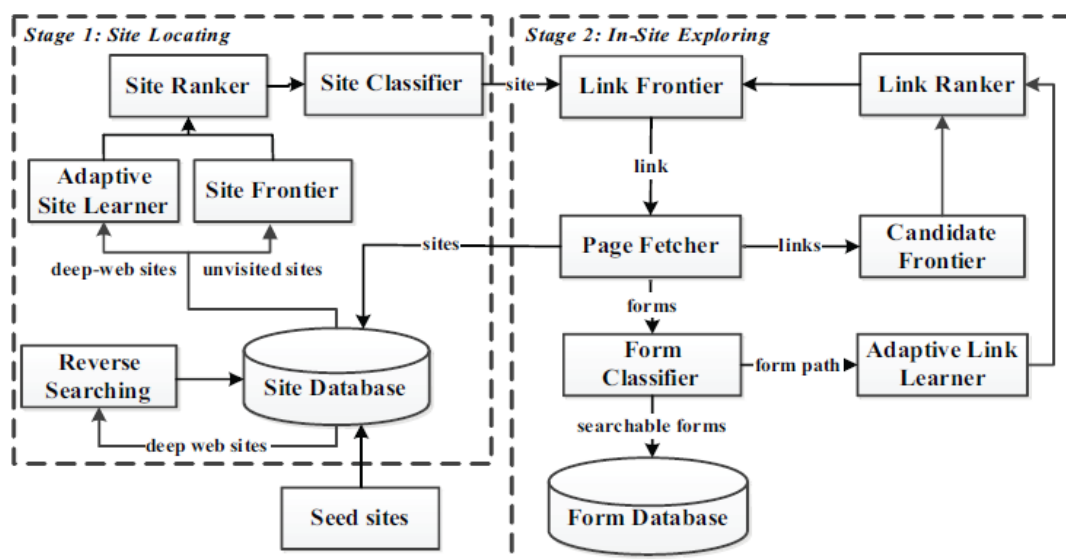
**A. Architecture:**



**Fig.1: System Architecture**

**Site Locating:**

The site locating stage finds relevant sites for a given topic, consisting of site colleting, site tanking, and site classification.

- **Site Collecting:**

The traditional crawler follows all newly found links. In contrast, our Smart Crawler strives to minimize the number of visited URLs, and at the same time maximizes the number of deep websites. To achieve these goals, using the links in

downloaded webpages is not enough. This is because a website usually contains a small number of links to other sites, even for some large sites. For instance, only 11 out of 259 links from webpages of aaronbooks.com pointing to other sites; amazon.com contains 54 such links out of a total of 500 links (many of them are different language versions, e.g., amazon.de). Thus, finding out-of-site links from visited web pages may not be enough for the Site Frontier. To address the above problem, we propose two crawling strategies, reverse searching and incremental two-level site prioritizing, to find more sites.

- **Reverse searching:**

The idea is to exploit existing search engines, such as Google, Baidu, and Bing etc., to find center pages of unvisited sites. This is possible because search engines rank web pages of a site and center pages tend to have high ranking values. Algorithm1 describes the process of reverse searching.

Pre-defined threshold. We randomly pick a known deep website or a seed site and use general search engine's facility to find center pages and other relevant sites, Such as Google's "link:" , Bing's "site:", Baidu's "domain:".

- **Incremental site prioritizing:**

 To make crawling process resumable and achieve broad coverage on websites, an incremental site prioritizing strategy is proposed. The idea is to record learned patterns of deep web sites and form paths for incremental crawling. First, the prior knowledge (information obtained during past crawling, such as deep websites, links with searchable forms, etc.) is used for initializing Site Ranker and Link Ranker.

- **Site Ranker**

Once the Site Frontier has enough sites, the challenge is how to select the most relevant one for crawling. In SmartCrawler, Site Ranker assigns a score for each unvisited site that corresponds to its relevance to the already discovered deep web sites.

- **Site Classifier:**

After ranking Site Classifier categorizes the site as topic relevant or irrelevant for a focused crawl, which is similar to page classifiers in FFC and ACHE. If a site is classified as topic relevant, a site crawling process is launched. Otherwise, the site is ignored and a new site is picked from the frontier. In SmartCrawler, we determine the topical relevance of a site based on the contents of its homepage. When a new site comes, the homepage content of the site is extracted and parsed by removing stop words and stemming. Then we construct a feature vector for the site and the resulting vector is fed into a Naive Bayes classifier to determine if the page is topic-relevant or not.

**In-Site Exploring:**

Once a site is regarded as topic relevant, in-site exploring is performed to find searchable forms. The goals are to quickly harvest searchable forms and to cover web directories of the site as much as possible. To achieve these goals, in-site exploring adopts two crawling strategies for high efficiency and coverage. Links within a site are prioritized with Link Ranker and Form Classifier classifies searchable forms.

- **Crawling Strategies:**

Two crawling strategies, stop-early and balanced link prioritizing, are proposed to improve crawling efficiency and coverage.

- **Stop-early:**

Previous work observed that 72% interfaces and 94% web databases are found within the depth of three. Thus, in-site searching is performed in breadth-first fashion to achieve broader coverage of web directories. Additionally, in-site searching employs the following stopping criteria to avoid unproductive crawling:

**SC1:** The maximum depth of crawling is reached.

**SC2:** The maximum crawling pages in each depth are reached.

**SC3:** A predefined number of forms found for each depth is reached.

**SC4:** If the crawler has visited a predefined number of pages without searchable forms in one depth, it goes to the next depth directly.

**SC5:** The crawler has fetched a predefined number of pages in total without searchable forms.

- **Balanced link prioritizing:**

The simple breadth-first visit of links is not efficient, whose results are in omission of highly relevant links and incomplete directories visit 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.

- **Link Ranker:**

Link Ranker prioritizes links so that Smart Crawler 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.

- **Form Classifier:**

Classifying forms aims to keep form focused crawling, which filters out non-searchable and irrelevant forms. For instance, an airfare search is often co-located with rental car and hotel reservation in travel sites. For a focused crawler, we need to remove off-topic search interfaces. SmartCrawler adopts the HIFI strategy to filter relevant searchable forms with a composition of simple classifiers. HIFI consists of two classifiers, a searchable form classifier (SFC) and a domain-specific form classifier (DSFC). SFC is a domain-independent classifier to filter out non-searchable forms by using the structure feature of forms. DSFC judges whether a form is topic relevant or not based on the text feature of the form that consists of domain-related terms. The strategy of partitioning the feature space allows selection of more effective learning algorithms for each feature subset. In our implementation, SFC uses decision tree based C4.5 algorithm and DSFC employs SVM.

- **Feature selection and ranking:**

SmartCrawler encounters a variety of web pages during a crawling process and the key to efficiently crawling and wide coverage is ranking different sites and prioritizing links within a site. This section first discusses the online feature construction of feature space and adaptive learning process of SmartCrawler, and then describes the ranking mechanism.

- **Online Construction of Feature Space:**

In SmartCrawler, patterns of links to relevant sites and searchable forms are learned online to build both site and link rankers. The ability of online learning is important for the crawler to avoid biases from initial training data and adapt to new patterns.

- **Adaptive Learning:**

SmartCrawler has an adaptive learning strategy that updates and leverages information collected successfully during crawling. As shown in Figure 1, both Site Ranker and Link Ranker are controlled by the site frequency measures the number of times a site appears in other sites. In particular, we consider the appearance in known deep sites to be more important than other sites.

- **Ranking Mechanism:**

- **Site Ranking:**

SmartCrawler 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. Site similarity measures the topic similarity between a new site and known deep web sites. Site frequency is the frequency of a site to appear in other sites, which indicates the popularity and authority of the site — a high frequency site is potentially more important. Because seed sites are carefully selected, relatively high scores are assigned to them.

- **Link Ranking:**

For prioritizing links of a site, the link similarity is computed similarly to the site similarity described above. The difference includes: 1) link prioritizing is based on the feature space of links with searchable forms (FSL); 2) for URL feature U, only path part is considered since all links have the same domain; and 3) the frequency of links is not considered in link ranking.

**B. Proposed Methodology:**

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 find many data sources for sparse domains. By ranking collected sites and by focusing the crawling on a topic, Smart Crawler achieves more accurate results.

## 3. RESULTS

**GUI:**

SmartCrawler: A Two-stage Crawler for Efficiently Harvesting Deep-Web Interfaces

| Home | Login | Registration | Admin |
| --- | --- | --- | --- |

**Login !**

| Email | : | | |
| Password | : | | |
| | | **Login** | |

**Output:**

SmartCrawler: A Two-stage Crawler for Efficiently Harvesting Deep-Web Interfaces

| Home | My Account | Edit Profile | Search Datasets | Logout |
| --- | --- | --- | --- | --- |

| Enter Keyword: | cake | Search |
| Enter Url: | Type your url here | Search |

SmartCrawler: A Two-stage Crawler for Efficiently Harvesting Deep-Web Interfaces

| Home | My Account | Edit Profile | Search Datasets | Logout |
| --- | --- | --- | --- | --- |

| File Name | Url | Rank | Action |
| --- | --- | --- | --- |
| cake1.txt | http://www.cake1.com | 4 | Download File |
| cake3.txt | http://www.cake3.com | 2 | Download File |
| cake2.txt | http://www.cake2.com | 1 | Download File |
| cake5.txt | http://www.cake5.com | 1 | Download File |
| cake4.txt | http://www.cake4.com | 0 | Download File |

click here To Advance Search

## 4. CONCLUSION

As profound web develops at a quick pace, there has been expanded enthusiasm for methods that assist proficiently with finding profound web interfaces. Nonetheless, because of the extensive volume of web assets and the dynamic way of profound web, accomplishing wide scope and high productivity is a testing issue. We propose a two-stage structure, in particular Smart Crawler, for effective gathering profound web interfaces. In the first stage, Smart Crawler performs site-based hunting down focus pages with the assistance of web indexes, abstaining from going by countless. To accomplish more exact results for an engaged slither, Smart Crawler positions sites to organize profoundly pertinent ones for a given point. In the second stage, Smart Crawler accomplishes quick in-site excavating so as to see most significant connections with a versatile connection positioning. To dispense with inclination on going by some exceedingly significant connections in shrouded web indexes, we outline a connection tree information structure to accomplish more extensive scope for a site. Our test results on an arrangement of delegate areas demonstrate the readiness and precision of our proposed crawler structure, which effectively recovers profound web interfaces from huge scale destinations and accomplishes higher harvest rates than different crawlers.

## ACKNOWLEDGMENTS

## 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 Raja Raman, 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 directory. http://www.clusty.com/, 2009