

DeepBot: A Focused Crawler for Accessing Hidden Web Content⁺

Manuel Álvarez, Juan Raposo, Alberto Pan*
Fidel CACHEDA, Fernando Bellas, Víctor Carneiro

Department of Information and Communications Technologies
University of A Coruña, Campus de Elviña s/n. 15071 A Coruña, Spain
{mad,jrs,apan,fidel,fbellas,viccar}@udc.es

Abstract. The crawler engines of today cannot reach most of the information contained in the Web. A great amount of valuable information is "hidden" behind the query forms of online databases, and/or is dynamically generated by technologies such as Javascript. This portion of the web is usually known as the Deep Web or the Hidden Web. We have built DeepBot, a prototype of hidden-web focused crawler able to access such content. DeepBot receives a set of domain definitions as an input, each one describing a specific data-collecting task and automatically identifies and learns to execute queries on the forms relevant to them. In this paper we describe the techniques employed for building DeepBot and report the experimental results obtained when testing it with several real world data collection tasks.

1 Introduction

A key component in the architecture of current Web search engines are the "crawler" programs used to automatically traverse the web, retrieving pages to build a searchable index of their content. Crawlers receive as input a set of "seed" pages and recursively obtain new ones by locating and traversing their outbound links.

Conventional web crawlers cannot reach to a very significant fraction of the web, which is usually called the "Hidden Web" or the "Deep Web".

Several works have studied and characterized the Hidden Web [4], [5]. They concluded that it is substantially larger than the publicly indexable web and that it usually contains data of higher quality and with a higher degree of structure.

The problem of crawling the "Hidden Web" can be divided into two challenges:

- *Crawling the "server-side" Hidden Web.* Many websites offer query forms to access the contents of an underlying database. Conventional crawlers cannot

⁺ This research was partially supported by the Spanish Ministry of Education and Science under project TSI2005-07730.

* Alberto Pan's work was partially supported by the "Ramón y Cajal" programme of the Spanish Ministry of Education and Science.

access these pages because they do not know how to execute queries on those forms.

- *Crawling the “client-side” Hidden Web.* Many websites use techniques such as client-side scripting languages and session maintenance mechanisms. Most conventional crawlers are unable to handle this kind of pages.

This paper overviews the architecture of DeepBot, a prototype system for crawling the Hidden Web, and describes in detail the techniques it uses for accessing the content behind web forms. The techniques used to deal with the client-side Deep Web were described in greater detail in [1].

The main features of DeepBot are:

- For accessing the “server-side” Deep Web, DeepBot can be provided with a set of *domain definitions*, each one describing a certain data-gathering task. DeepBot automatically detects forms relevant to the defined tasks and executes a set of pre-defined queries on them.
- DeepBot’s crawling processes are based on automated “mini web browsers”, built by using browser APIs (our current implementation is based on Microsoft Internet Explorer). This enables our system to deal with client-side scripting code, session mechanisms, and other complexities related with the client-side Hidden Web.

The paper is organized as follows. Section 2 overviews the architecture of DeepBot and the main components that participate in accessing the server-side Hidden Web. Section 3 describes the domain definitions used to specify a data collection task. Section 4 describes how DeepBot detects query forms relevant to a certain task and how it learns to execute queries on them. Section 5 describes our experiments with the system. Section 6 discusses related work and section 7 concludes the paper.

2 Architecture

The architecture of the system is shown in Fig. 1. As well as in conventional crawlers, the functioning of DeepBot is based on a shared list of *routes* (pointers to documents), which will be accessed by a certain number of concurrent crawling processes, distributed into several machines.

The main singularities of our approach are:

- In conventional crawlers, routes are just URLs. Thus, they have problems with sources using session mechanisms. Our system stores, with each route, a session object containing all the required information (cookies, etc.) to restore the execution environment in which the crawling process was running in the moment of adding the route to the master list.
- Conventional engines implement crawling processes by using http clients. Instead, our system uses lightweight automated *mini web browsers* (built by using the APIs of most popular browsers) as execution environment for automated navigation. These *mini web browsers* access to pages by generating actions on a web browser interface, in the same way a human user would generate them when browsing. For

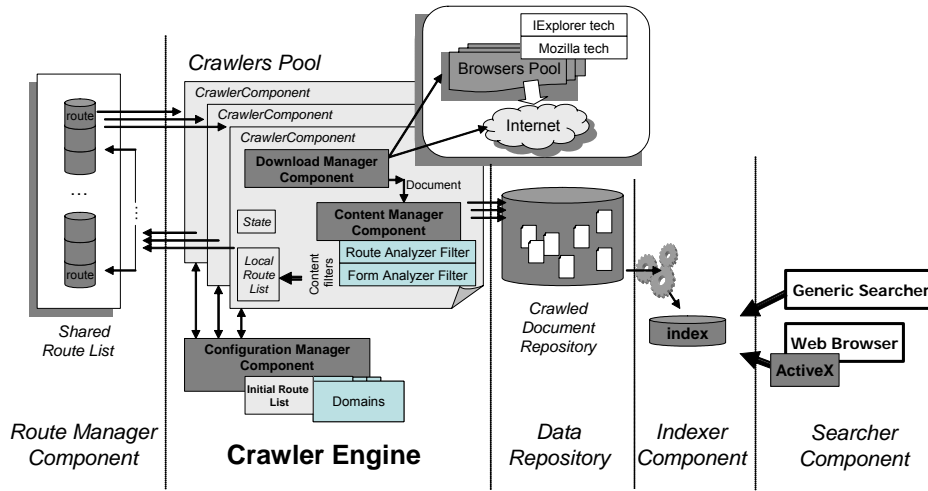


Fig. 1. Crawler Architecture

specifying a navigation sequence in the automated mini-browsers, we have created NSEQL [16], a language which allows representing a sequence as the list of interface events a user would need to produce on the web browser in order to reach the desired page.

- When the system reaches a new page, in addition of using its anchors to generate new routes, it also examines each HTML form and ranks its relevance with respect to a set of pre-configured *domain definitions* (remember that each domain definition describes a specific data-collection task). If the system finds that the form is relevant, it is used to execute a set of queries defined by the domain, thus reaching to new pages.

The architecture also includes components for indexing and searching the crawled contents, using state of the art algorithms (our current implementation is based on Apache Lucene). The NSEQL sequence needed to access each document is also stored. This sequence is used by the *ActiveX for automatic navigation Component*, which receives as a parameter a NSEQL program, downloads itself into the user browser and makes it execute the given sequence. This is used to access the documents returned as result of a search against the index, when they cannot be directly accessed in the source by using its URL, due to session issues.

3 Domain Definitions

In this section, we describe the domain definitions used to define a data-collection task. A *domain definition* is composed of the following elements:

- A set of *attributes* $A = \{a_1, a_2, \dots, a_n\}$. Each attribute a_i has associated a *name*, a set of aliases $\{a_{i_alias_1}, \dots, a_{i_alias_k}\}$, and a specificity index s_i .

Domain "Books"			
Attributes: $A = \{a_1, a_2, a_3, a_4, a_5\}$			
Attribute	Name	Aliases	s_i (specificity index)
a_1	TITLE	'entitle', 'title of the book'	0.6
a_2	AUTHOR	'writer', 'written by'	0.7
a_3	ISBN		1
a_4	FORMAT	'binding type'	0.25
a_5	PRICE	'cost of book'	0.05
Queries: $Q = \{q_1, q_2, q_3, q_4, q_5, q_6, q_7, q_8\}$			
$q_1 = \{(\text{TITLE}, \text{'java'}), (\text{FORMAT}, \text{'hardcover'})\}$			
$q_2 = \{(\text{TITLE}, \text{'xml'}), (\text{AUTHOR}, \text{'Priscilla Walmsley'})\}$			
$q_3 = \{(\text{TITLE}, \text{'jee'})\}$			
$q_4 = \{(\text{TITLE}, \text{'concurrent programming'})\}$			
$q_5 = \{(\text{TITLE}, \text{'ejb3'})\}$			
$q_6 = \{(\text{TITLE}, \text{'java server faces'})\}$			
$q_7 = \{(\text{AUTHOR}, \text{'Herbert Schildt'})\}$			
$q_8 = \{(\text{TITLE}, \text{'web services'})\}$			
Relevance threshold: $\mu = 0.9$			

Fig. 2. Example query form and visual distances and angles for field f_1

- A set of *queries* $Q = \{q_1, q_2, \dots, q_m\}$ we want to execute on the discovered relevant forms. Each query q_j is a list of pairs (*attribute*, *value*), where *attribute* is an attribute of the domain and *value* is a string (it can be empty).
- A relevance threshold denoted as μ .

An *attribute* represents a field that may appear in the query forms that are relevant to the data-collection task.

The *aliases* represent alternative labels that may identify the attribute in a query form. For instance, the attribute AUTHOR, from a domain used for collecting data about books, could have aliases such as “*writer*” or “*written by*”. It is important to notice that the study in [5] concluded that the aggregate schema vocabulary of web forms in the same domain tends to converge at a relatively small size. They also detected a Zipf-like distribution of attribute frequencies (thus, a small set of “dominant” attributes are much more frequent in the forms of the domain than the rest of attributes). This supports the feasibility of creating effective domain specifications in an easy and fast way: exploring a few sources in the domain is usually enough to find the most important attributes and aliases.

For each attribute, the domain also includes a *specificity index*. The specificity index (denoted s_i) of an attribute a_i is a number between 0 and 1 indicating how probable is that a query form containing such attribute is actually relevant to the domain. For instance, in an example domain for collecting book data, the attribute ISBN would have a very high value (e.g. 0.95), since a query form allowing queries for the ISBN attribute is almost certainly a form allowing to search books; the PRICE attribute would have a low value such as 0.05, since a query form containing it could be related to any kind of product.

Finally, the domain also includes a *relevance threshold* μ . The specificity indexes and the threshold will be used to determine if a given form is relevant to a domain.

Fig. 2 shows an example domain definition for the task of collecting pages containing data about books on the subject of Java and XML programming. The relevance threshold for this domain is set to 0.9.

4 Processing Forms with the Form Analyzer

In this section, we describe how the crawler processes each found form. The performed steps are:

- For every domain, the system tries to match its attributes with the fields of the form, using visual distance and text similarity heuristics (see subsection 4.1).
- By using the output of the previous step, the system determines if the form is relevant with respect to the domain (described in subsection 4.2).
- If the form is relevant, the crawler uses it to execute the queries defined in the domain. For each query, we obtain a new route to add to the list of routes. The new route will be dealt with as any other route fetched by the crawler (subsection 4.3).

4.1 Associating Form Fields and Domain Attributes

Given a form f located in a certain HTML page and a domain d describing a data-collecting task, our goal at this stage is to determine whether f allows executing queries for the attributes of the domain d or not. Our method consists of two steps:

1. Determining which texts are associated with each field of the form. This step is based on heuristics using visual distance measures between the form fields and the texts surrounding them.
2. Trying to relate the fields of f with the attributes of d . The system performs this step by obtaining text similarity measures between the texts associated with each form field and the texts associated with each attribute in the domain definition d .

Measuring visual distances. At this step, we consider the texts in the page and compute their visual distance with respect to each field of the form f . The visual distance between a text element t and a form field f is computed as follows:

1. The browser APIs are used to obtain the coordinates of a rectangle enclosing f and a rectangle enclosing t . If t is into an HTML table cell, and it is the unique text inside, then the coordinates of the table cell rectangle are assigned to t .
2. We obtain the minimum distance between both rectangles. Distances are not computed in pixels but in more coarse-grained units (we use cells of the approximated visual size of one character).
3. We also obtain the angle of the shortest line joining both rectangles. The angle is approximated to the nearest multiple of $\pi/4$.

Fig. 3a shows one example query form corresponding to an Internet bookshop. We show the distance and angles obtained for some of its texts and fields.

Associating texts and form fields. For each form field, our goal is to obtain the texts “semantically linked” with it in the page. For instance, in the Fig. 3a the strings semantically linked to the first field are “*Book Title*” and “(example: ‘*Thinking in Java*’)”. For pre-selecting the “best texts” for a field f , we apply the following steps:

1. We add all the texts having the shortest distance d with respect to f to the list.

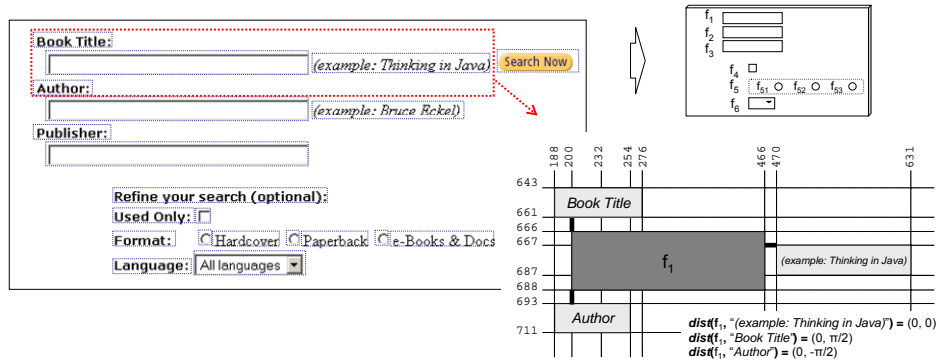


Fig. 3a. Example query form and visual distances and angles for field f_1

2. Those texts having a distance lesser than $k \cdot d$ with respect to f are added to the list ordered by distance (k is a configurable factor usually set to 5). This step discards those texts that are significantly further from the field.
3. Texts with the same distance are ordered according to its angle. The preference order for angles privileges texts aligned with the fields (that is, angle multiple of $\pi/2$); it also privileges left with respect to right and top with respect to bottom, because they are the preferred positions for labels in forms.

As output of the previous step we have an ordered list of texts, which are probably associated to each form field. Then we post-process the lists as follows:

1. We ensure that a given text is only present in the list of one field. The rationale for this is that at the following stage of the form ranking process (which consists in matching form fields and “searchable” attributes), we will need to associate unambiguously a certain text with a given form field.
2. We ensure that each field has at least one associated text. The rationale for this is that, in real pages, a given form field always has some associated text to allow the user to identify its function. For instance, if the list of a field f_1 contained the texts t_1 and t_2 (in that order), and the list of a field f_2 only contained the text t_1 , then we would choose to remove t_1 from the list of f_1 , since removing it from the list of f_2 would leave the field with an empty list. Note that this would be done even if t_1 were actually closer to f_1 than to f_2 .

Fig. 3b shows the process for the example form of Fig. 3a. For each field¹ of the form, we show the ordered list of texts obtained by applying the visual distance and angle heuristics. The texts remaining in the lists after the post-processing steps are boldfaced in the figure.

Associating form fields and domain attributes. At this step we try to detect the form fields which correspond to attributes of the target domain. We distinguish between two kinds of fields:

- *Bounded fields.* We term as bounded those fields offering a finite list of possible query values, such as *select-option* fields, *checkbox* fields or *radio buttons*.

¹ Note how the system models the FORMAT ‘checkbox’ field as a field with three subfields. f_5 refers to the whole set of *checkboxes* while f_{51} , f_{52} and f_{53} refer to individual *checkboxes*.

Fields	Texts	(dist, θ)	Fields	Texts	(dist, θ)	Fields	Texts	(dist, θ)
f_1	√ (example: Thinking in Java)	(0, 0)	f_4	√ Used Only:	(0, π)	f_{52}	Hardcover	(0, π)
	√ Book Title:	(0, $\pi/2$)		√ Refine your search (optional):	(0, $\pi/2$)		√ Paperback	(0, 0)
	Author:	(0, $-\pi/2$)		Hardcover	(0, $-\pi/2$)	f_{53}	Refine your search (optional):	(1, $\pi/2$)
	(example: Bruce Eckel)	(1, $-\pi/2$)		Format:	(1, π)		Paperback	(0, π)
f_2	√ (example: Bruce Eckel)	(0, 0)	f_5	e-Books & Docs	(0, 0)	f_6	√ e-Books & Docs	(0, 0)
	√ Author:	(0, $\pi/2$)		Used Only:	(0, $3\pi/4$)		Refine your search (optional):	(2, $3\pi/4$)
	Publisher:	(0, $-\pi/2$)		Language	(0, $-3\pi/4$)		Language:	(0, π)
	(example: Thinking in Java)	(2, $\pi/2$)		√ Format:	(1, π)		Hardcover	(0, $\pi/2$)
f_3	√ Book Title:	(3, $\pi/2$)	f_{51}	√ Refine your search (optional):	(1, $\pi/2$)		Paperback	(0, $\pi/2$)
	√ Publisher:	(0, $\pi/2$)		√ Hardcover	(0, 0)		Format:	(1, π)
	Refine your search (optional):	(1, $-\pi/2$)		Used Only:	(0, $3\pi/4$)		Used Only:	(1, $\pi/2$)
	(example: Bruce Eckel)	(2, $\pi/2$)		Language:	(0, $-3\pi/4$)		Refine your search (optional):	(3, $\pi/2$)
	Author:	(3, $\pi/2$)		Format:	(1, π)			
	Used Only:	(3, $-\pi/2$)		Refine your search (optional):	(1, $\pi/2$)			
	Format:	(4, $-\pi/2$)						
	Hardcover	(4, $-\pi/2$)						
	Paperback	(4, $-\pi/2$)						

f_1 [(example: Thinking in Java)] [Book Title:] [Author:]; f_2 [(example: Bruce Eckel)] [Author:]; f_3 [Publisher:];
 f_4 [Used Only:] [Refine your search (optional):];
 f_5 [Format:]; f_{51} [Hardcover]; f_{52} [Paperback]; f_{53} [e-Books & Docs]; f_6 [Language]

Fig. 3b. Texts associated to each field in the form of Fig. 3a

- *Unbounded fields*. We term as unbounded those fields whose query values are not limited, such as *text boxes*.

The basic idea to rank the “similarity” between a field f and an attribute a is to measure the textual similarity between the texts associated with f in the page (obtained as shown in the previous step) and the texts associated with a in the domain (the attribute name and the aliases). When the field is *bounded*, the system also takes into account the text similarities between the possible values of f in the page² and the query input values specified for a in the domain queries. Text similarity measures are obtained using a method proposed in [7] that combines TFIDF and the Jaro-Winkler edit-distance algorithm. Before computing similarities, the texts are normalized by removing *stopwords* and non alphanumeric characters. The words in the text are ordered alphabetically. Stemming techniques are not used because word suffixes are very useful to differentiate aliases (e.g. “*Published*” and “*Publisher*”).

We use the following method to rank the *similarity* between a field f and a searchable attribute a :

1. We indicate the texts describing a (that is, the attribute name and its aliases specified in the domain) as the set $\{a_alias_1, \dots, a_alias_n\}$. We obtain the texts associated with f in the page using the methods seen in the previous section; we notate them as $\{f_text_1, \dots, f_text_m\}$.
2. If f is a *bounded* field, then we also obtain its possible values from the page $\{f_value_1, \dots, f_value_p\}$. We also examine the queries from the domain to obtain the set of input values used in them for attribute a $\{a_value_1, \dots, a_value_q\}$.
3. Now we compute the text similarity between each pair $(a_alias_i, f_text_j)_{i=1..n, j=1..m}$. Then we obtain sim_1 as in (1), i.e. the maximum of the obtained text similarities.

² Obtaining these values is a trivial step for *select-option* tags, since their possible values appear in the HTML code enclosed in *option* tags. For *checkbox* and *radio* tags we apply visual distance techniques similar to the ones previously discussed.

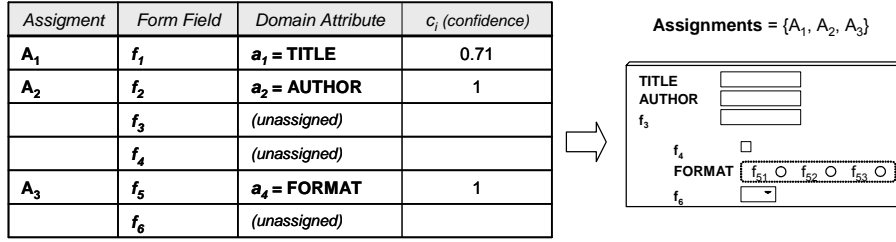


Fig. 4. Assignments obtained for the form in Fig. 3a, using the domain definition shown in Fig. 2

$$sim_1 = \max \left\{ textSim(a_alias_i, f_text_j) \right\}_{i=1..n, j=1..m} \quad (1)$$

4. If f is bounded, then we also compute the text similarity between each pair $(f_value_k, a_value_r)_{k=1..p, r=1..q}$. Then we obtain sim_2 as in (2), i.e. the mean of the maximum similarities obtained for each $a_value_{r, r=1..q}$.

$$sim_2 = \sum_{r=1..q} \left(\max \left\{ textSim(a_value_r, f_value_k) \right\}_{k=1..p} \right) / q \quad (2)$$

5. If f is bounded, then the similarity between a and f is set to $\max \{sim_1, sim_2\}$.
If f is unbounded, then the similarity is set to sim_1 .

As result of applying the previous steps, we obtain a table with the estimated similarities between each form field and each attribute. Then we proceed as follows:

- The pairs from the table that do not reach a minimum similarity *threshold* are discarded.
- If the table does not contain two entries for the same attribute, the process finishes and the table contains the valid associations between form fields and attributes. If the table contains more than one entry for the same attribute, we choose for each attribute the entry with a higher similarity but trying to guarantee that no field with an entry above the threshold is left unassigned.

The output of this stage is a set of assignments between form fields and domain attributes. Each of these assignments has a certain *confidence*, which the system sets to the similarity obtained between the field and the attribute.

Fig. 4 shows the assignments obtained for the form in Fig. 3a, using the domain definition of Book Shopping shown in Fig. 2.

4.2 Determining the Relevance of a Form to a Domain

The output of the previous stage is a set of assignments $\{A_1, \dots, A_k\}$ between form fields and domain attributes. Each of these assignments has a certain *confidence*, expressed as a number between 0 and 1. We indicate the confidence of assignment A_i as c_i .

The method we use to determine if a form is relevant to a domain consists of adding the confidences of each assignment, pondered by the *specificity index* of the attribute involved in it, and checking if the sum exceeds the *relevance threshold* μ . That is, the system checks if the inequality in (3) is verified.

"Books Shopping"		
Attributes		
Attribute Name	Aliases	s_i (specificity index)
TITLE	'title of book'	0.6
AUTHOR	'author's name'	0.7
PUBLISHER		0.8
ISBN		0.95
PUBDATE	'publication date'	0.7
SUBJECT	'section', 'category', 'department', 'subject Category'	0.05
FORMAT	'binding type'	0.25
PRICE		0.05
Relevance threshold: $\mu = 0.9$		

"Music Shopping"		
Attributes		
Attribute Name	Aliases	s_i (specificity index)
ARTIST	'artist name', 'composer/author/artist'	0.6
SONG	'soundtrack title', 'song title'	0.95
ALBUM	'album title'	0.95
LABEL	'vendor'	0.8
GENRE	'style'	0.05
FORMAT	'media type', 'product type', 'item types'	0.25
PRICE		0.05
Relevance threshold: $\mu = 0.9$		

"Movies Shopping"		
Attributes		
Attribute Name	Aliases	s_i (specificity index)
TITLE	'movie title'	0.6
LEGEND		0.7
STARRING	'star', 'actor', 'cast', 'featuring (cast/crew)', 'cast name', 'artistic'	0.7
DIRECTOR		0.7
PRODUCER		0.7
EDITOR		0.7
SOUND	'music'	0.7
FORMAT	'media'	0.05
GENRE	'movie type', 'category'	0.05
PRICE		0.05
Relevance threshold: $\mu = 0.9$		

Fig. 5. Domain definitions: Books, Music and Movies

$$\sum_{i=1..k} c_i s_i > \mu \quad (3)$$

For instance, considering the domain definition shown in Fig. 2, and the assignments shown in Fig. 4, we would obtain (4).

$$0.71 \cdot 0.6 + 1 \cdot 0.7 + 1 \cdot 0.25 = 1.376 > \mu = 0.9 \quad (4)$$

4.3 Executing Queries

Once the system determines that a form is relevant to a certain domain d , a new route must be added for each query specified in d . Executing a query in the form involves filling in the form according to the query and submitting the form.

The first task can be easily done from the assignments which associate form fields and domain attributes.

The second task has its own complications. Although the lightweight mini-browsers the system uses as crawling processes may directly issue a SUBMIT event on the form once it has been filled in, this simple strategy does not work in some websites. This is due to the frequent use of client-side scripting languages to manage form submission. To overcome these difficulties, the system proceeds as follows:

1. The system searches for *input* elements in the form of the types *submit*, *image* or *button* (in that order). Each found element is used to try to submit the form by generating a *click* event on it. After each try, the system checks if the event caused a new navigation in the browser. If it was not the case, it tries the next element.
2. If the previous step is unsuccessful (typically because the searched types of *input* elements do not exist), the system concludes that the way used to submit the form is clicking on an anchor with some associated client-scripting code (typically Javascript). Therefore, the system looks for anchors located visually close to the form and having associated some client-side script in either the *href* or the *onClick* attributes. The anchors obtained are ordered according to its visual proximity to the form and to the text similarity between their associated texts and a set of pre-defined texts commonly used to indicate form submission

- (e.g. ‘search’, ‘go’, ‘submit’,...). The system tries to generate a *click* event on the anchors in the list and checks if the event caused a new navigation in the browser. If it is not the case, it tries the next element.
3. If all the previous steps fail, the system generates a SUBMIT event on the form.

5 Experience

To evaluate the performance of our approach, we tested it on three different domains: Books Shopping, Music Shopping and Movies Shopping websites.

The process for creating the domain definitions was the following: for each domain, we manually explored 10 sites at random, from the respective Yahoo Directory³ category and used them to define the attributes and aliases. The specificity indexes and the relevance threshold were also manually chosen from our experience visiting these sites. The resulting domain definitions are shown in Fig. 5.

Once the domains were created, we used DeepBot to crawl 20 websites of the respective Yahoo Directory category. The websites visited by DeepBot for each domain are shown in Table 1. The websites used to define the attributes and aliases are grouped in a dataset named *Training*, while the remaining sites are grouped in a dataset named *Advanced*.

To check the accuracy of the results obtained, we manually analyzed the websites and compared the results with those obtained by DeepBot. We measured the results at each stage of the process: associating texts with form fields, associating form fields with domain attributes, establishing the relevance of a form to a domain, and executing the queries on the relevant forms.

To quantify the results, we used standard Information Retrieval metrics: precision, recall and F_1 -measure. The metrics defined to measure the performance of DeepBot, make use of the following variables:

- $TextFieldA_{DeepBot}$: set of the associations between texts and form fields discovered by DeepBot.
- $TextFieldA_{Real}$: set of the associations between texts and form fields discovered by the manual analysis.
- $FieldAttributeA_{DeepBot}$: set of the associations between form fields and domain attributes discovered by DeepBot.
- $FieldAttributeA_{Real}$: set of the associations between form fields and domain attributes discovered by the manual analysis.
- $FormDomainA_{DeepBot}$: set of the associations between forms and domains discovered by DeepBot.
- $FormDomainA_{Real}$: set of the associations between forms and domains discovered by manual analysis.
- $SubmittedForms_{DeepBot}$: set of forms successfully submitted by DeepBot.

³ <http://dir.yahoo.com>

Table 1. Training and advanced datasets for Books, Music and Movies Shopping domains

Domain \ DataSet	Training (S_1)	Advanced (S_2)
Books Shopping	Amazon.com - http://www.amazon.com Barnes&Noble - http://www.barnesandnoble.com Powell’s Books - http://www.powells.com eCampus.com - http://www.ecampus.com Dymocks Booksellers - http://www.dymocks.com.au Tattered Cover Bookstore - http://www.tatteredcover.com BookFinder4U.com - http://www.bookfinder4u.com Cody’s Books - http://www.codysbooks.com Daedalus Books&Music - http://www.daedalusbooks.com Bolen - http://www.bolen.bc.ca	The Book Pl@ce - http://www.thebookplace.com Blackwell’s Bookshop - http://bookshop.blackwell.co.uk The American Book Center - http://www.abc.nl Oxbow Books - http://www.oxbowbooks.com Strand Book Store - http://www.strandbooks.com Globe Pequot Press - http://www.globepequot.com Northshire Bookstore - http://www.northshire.com Green Apple Books&Music- http://www.greenapplebooks.com Thomson Gale - http://www.gale.com The Scholar’s Bookshelf - http://www.scholarsbookshelf.com
Music Shopping	Amazon.com Music - http://www.amazon.com Barnes&Noble Music - http://www.barnesandnoble.com Tower Records - http://www.towerrecords.com Rough Trade - http://www.roughtrade.com Sam Goody - http://www.samgoody.com Schott Musik International – http://www.schott-music.com A&B Sound - http://catalog2.absound.ca Collectors’ Choice Music - http://www.ccmusic.com CD Quest - http://www.cdquest.com AudibleFaith – http://www.audiblefaith.com	Cyber Music Surplus - http://www.cybermusic surplus.com ClassicTrax.co.uk - http://www.classictrax.ltd.uk MyMusic.com - http://www.mymusic.com Mojo Sounds - http://www.mojosounds.com Looney Tunes - http://www.looneytunes cds.com Buy Cd - http://www.buycd.com Record Exchange - http://www.buymusic here.net Musica Obscura - http://www.musicaobscura.com ProMusicFind.com - http://www.promusic find.com Ladyslipper Music - http://www.ladyslipper.org
Movies Shopping	Amazon.com DVD – http://www.amazon.com The Spinning Image - http://www.thespinningimage.co.uk/ Saregama - http://hamaracd.com Docuseek - http://www.docuseek.com BollyVista.com - http://www.bollyvista.com Movie Tickets.com - http://www.movietickets.com Yahoo! Movies - http://movies.yahoo.com Half.com - http://www.half.ebay.com GlobeAndMail.com - http://www.theglobeandmail.com Kansas.com - http://ae.kansas.com	Fandango - http://www.fandango.com Sensasian.com - http://sensasian.com Blockbuster.com - http://www.blockbuster.com Video*ezy - http://www.videoezy.com.au Intelliflix - http://www.intelliflix.com IGN.COM - http://search.ign.com Bestvideobuys - http://www.bestwebbuys.com BrickFilms - http://www.brickfilms.com Human Rights Film Directory - http://db.lib.washington.edu Blockbuster UK - http://www.blockbuster.co.uk

We defined the following metrics:

- Metrics for associating texts and form fields in (5).

$$\begin{aligned}
 Precision_{TextFieldA} &:= \frac{|TextFieldA_{DeepBot} \cap TextFieldA_{Real}|}{|TextFieldA_{DeepBot}|} \\
 Recall_{TextFieldA} &:= \frac{|TextFieldA_{DeepBot} \cap TextFieldA_{Real}|}{|TextFieldA_{Real}|} \\
 F_1 - measure_{TextFieldA} &:= 2 \times Precision_{TextFieldA} \times Recall_{TextFieldA} / (Precision_{TextFieldA} + Recall_{TextFieldA})
 \end{aligned} \tag{5}$$

- Metrics for associating form fields and domain attributes in (6).

$$\begin{aligned}
Precision_{FieldAttributeA} &:= |\text{FieldAttributeA}_{DeepBot} \cap \text{FieldAttributeA}_{Real}| / |\text{FieldAttributeA}_{DeepBot}| \\
Recall_{FieldAttributeA} &:= |\text{FieldAttributeA}_{DeepBot} \cap \text{FieldAttributeA}_{Real}| / |\text{FieldAttributeA}_{Real}| \\
F_1 - measure_{FieldAttributeA} &:= 2 \times Precision_{FieldAttributeA} \times Recall_{FieldAttributeA} / (Precision_{FieldAttributeA} + Recall_{FieldAttributeA})
\end{aligned} \tag{6}$$

- Metrics for Global associations between forms and domains in (7).

$$\begin{aligned}
Precision_{FormDomainA} &:= |\text{FormDomainA}_{DeepBot} \cap \text{FormDomainA}_{Real}| / |\text{FormDomainA}_{DeepBot}| \\
Recall_{FormDomainA} &:= |\text{FormDomainA}_{DeepBot} \cap \text{FormDomainA}_{Real}| / |\text{FormDomainA}_{Real}| \\
F_1 - measure_{FormDomainA} &:= 2 \times Precision_{FormDomainA} \times Recall_{FormDomainA} / (Precision_{FormDomainA} + Recall_{FormDomainA}) \\
Precision_{SubmittedForms} &:= |\text{SubmittedForms}_{DeepBot}| / |\text{FormDomainA}_{DeepBot} \cap \text{FormDomainA}_{Real}|
\end{aligned} \tag{7}$$

5.1 Experimental Results

Table 2 summarizes the obtained experimental results. For each domain, it shows the values obtained for all the metrics in the *Training* dataset (S_1 , the sites used to define the domains), the *Advanced* dataset (S_2 , the remaining sites) and in the *Global* dataset ($S_1 + S_2$, *Training* + *Advanced*).

It is important to notice that, in order to calculate the metrics for form-domain and field-attribute associations, “quick search” and authentication forms have not been considered. The results include only multi-field forms of the kind usually employed for “advanced search” forms. In addition, the results for the field-attribute associations have been measured independently of the results of the previous stage (text-field associations).

The obtained results are quite promising: all the metrics show high values and some of them even reach 100%. Now we discuss the reasons behind the mistakes committed by DeepBot at each stage.

Recall in associating forms and domains reached 100% in every case but in the *Advanced* dataset of the Music and Movies domains (which reached 95%). In the music domain, the reason was that the *ProMusicFind* source used an alias for the ARTIST attribute which did not match with any of the aliases defined in the domain. In addition, the form only had two fields so, even though the system correctly assigned the other one to a domain attribute (“Album Title”), it was not enough to exceed the relevance threshold. In the movies domain, the query form from source *IGN.COM* only had two searchable fields (title, genre) matching with attributes in our domain definition. Although the system correctly matched both, it was not enough to reach the threshold.

The precision and recall values obtained for the associations between texts and form fields exceeded 80% except in the *Advanced* dataset of the Books domain (0.73 precision and 0.79 recall). The majority of the errors in this dataset came from a single source (*Blackwell's Bookshop*). If we did not have into account this source, the metrics would take values similar to those reached by the other ones.

Table 2. Experimental results

	Books Shopping			Music Shopping			Movies Shopping		
	S ₁	S ₂	S ₁ +S ₂	S ₁	S ₂	S ₁ +S ₂	S ₁	S ₂	S ₁ +S ₂
Submitted Forms									
Precision	13/13 1.00	11/11 1.00	24/24 1.00	10/10 1.00	9/9 1.00	19/19 1.00	12/12 1.00	9/9 1.00	21/21 1.00
Form-Domain Associations									
Precision	13/13 1.00	11/11 1.00	24/24 1.00	10/10 1.00	9/9 1.00	19/19 1.00	12/12 1.00	9/9 1.00	20/20 1.00
Recall	13/13 1.00	11/11 1.00	24/24 1.00	10/10 1.00	9/10 0.90	19/20 0.95	12/12 1.00	9/10 0.90	21/22 0.95
F₁-measure	1.00	1.00	1.00	1.00	0.95	0.97	1.00	0.95	0.97
Field-Attribute Associations									
Precision	54/55 0.98	50/50 1.00	104/105 0.99	37/37 1.00	31/33 0.94	68/70 0.97	45/46 0.98	33/33 1.00	78/79 0.99
Recall	54/54 1.00	50/53 0.94	104/107 0.97	37/37 1.00	31/37 0.84	68/74 0.92	45/45 1.00	33/35 0.94	78/80 0.98
F₁-measure	0.99	0.97	0.98	1.00	0.89	0.94	0.99	0.97	0.98
Text-Field Associations									
Precision	129/142 0.91	101/137 0.73	230/279 0.82	93/110 0.83	107/132 0.81	199/242 0.82	154/179 0.86	163/184 0.89	317/363 0.87
Recall	129/132 0.98	101/127 0.79	230/259 0.88	92/94 0.98	107/109 0.98	199/203 0.98	154/168 0.92	163/181 0.90	317/349 0.91
F₁-measure	0.94	0.76	0.85	0.90	0.89	0.89	0.89	0.89	0.89

The failures at this stage came mainly from *bounded* fields that did not have any globally associated text in the form (the form only included the texts corresponding to its values). That is contrary to one of our heuristics, which assumed that every form field should have at least one associated text to “explain” the function of the field to the user. It actually turns out that the values of the bounded field may be auto-explicative and, in some rare cases, the form creators choose to not include an additional text label.

For instance, from a list of *radio buttons* including the options “*hardcover*” and “*paperback*”, the user may infer the buttons are used to express a query condition on the format attribute, even when there is not additional text label indicating it.

Finally, *Recall* and *Precision* also reach high values (> 90% except in one case) in the associations between form fields and domain attributes. The mistakes at this stage typically occurred because the domain did not include the alias used in the form for some attribute.

Tables 3, 4 and 5 show the obtained experimental results for each website.

Table 3. Experimental results for each site in *Books Shopping* domain.

Datasets / Metrics	Form-Domain Associations			Field-Attribute Associations		Text-Field Associations	
	Precision	Recall	Precision	Precision	Recall	Precision	Recall
			SubmittedForms				
Amazon.com	1/1 = 1.00	1/1 = 1.00	1/1 = 1.00	7/7 = 1.00	7/7 = 1.00	17/19 = 0.89	17/17 = 1.00
Barnes&Noble	2/2 = 1.00	2/2 = 1.00	2/2 = 1.00	6/7 = 0.86	6/6 = 1.00	9/11 = 0.82	9/9 = 1.00
Powell's Books	1/1 = 1.00	1/1 = 1.00	1/1 = 1.00	6/6 = 1.00	6/6 = 1.00	16/17 = 0.94	16/16 = 1.00
eCampus.com	1/1 = 1.00	1/1 = 1.00	1/1 = 1.00	4/4 = 1.00	4/4 = 1.00	7/7 = 1.00	7/7 = 1.00
Dymocks Booksellers	1/1 = 1.00	1/1 = 1.00	1/1 = 1.00	4/4 = 1.00	4/4 = 1.00	7/8 = 0.88	7/7 = 1.00
Tattered Cover Bookstore	2/2 = 1.00	2/2 = 1.00	1/1 = 1.00	5/5 = 1.00	5/5 = 1.00	7/8 = 0.88	7/7 = 1.00
BookFinder4U.com	1/1 = 1.00	1/1 = 1.00	1/1 = 1.00	6/6 = 1.00	6/6 = 1.00	24/27 = 0.89	24/25 = 0.92
Cody's Books	2/2 = 1.00	2/2 = 1.00	2/2 = 1.00	5/5 = 1.00	5/5 = 1.00	7/8 = 0.88	7/7 = 1.00
Daedalus Books&Music	1/1 = 1.00	1/1 = 1.00	1/1 = 1.00	7/7 = 1.00	7/7 = 1.00	30/32 = 0.94	30/32 = 0.94
Bolen	1/1 = 1.00	1/1 = 1.00	1/1 = 1.00	4/4 = 1.00	4/4 = 1.00	5/5 = 1.00	5/5 = 1.00
The Book Pl@ce	1/1 = 1.00	1/1 = 1.00	1/1 = 1.00	4/4 = 1.00	4/4 = 1.00	8/9 = 0.89	8/9 = 0.89
Blackwell's Bookshop	1/1 = 1.00	1/1 = 1.00	1/1 = 1.00	7/7 = 1.00	7/7 = 1.00	12/33 = 0.36	12/28 = 0.43
The American Book Center	1/1 = 1.00	1/1 = 1.00	1/1 = 1.00	5/5 = 1.00	5/6 = 0.83	13/13 = 1.00	13/13 = 1.00
Oxbow Books	2/2 = 1.00	2/2 = 1.00	2/2 = 1.00	6/6 = 1.00	6/6 = 1.00	12/15 = 0.80	12/12 = 1.00
Strand Book Store	1/1 = 1.00	1/1 = 1.00	1/1 = 1.00	7/7 = 1.00	7/8 = 0.88	11/13 = 0.85	11/12 = 0.92
Globe Pequot Press	1/1 = 1.00	1/1 = 1.00	1/1 = 1.00	5/5 = 1.00	5/5 = 1.00	6/6 = 1.00	6/6 = 1.00
Northshire Bookstore	1/1 = 1.00	1/1 = 1.00	1/1 = 1.00	6/6 = 1.00	6/6 = 1.00	13/14 = 0.93	13/14 = 0.93
Green Apple Books&Music	1/1 = 1.00	1/1 = 1.00	1/1 = 1.00	3/3 = 1.00	3/4 = 0.75	6/9 = 0.67	6/9 = 0.67
Thomson Gale	1/1 = 1.00	1/1 = 1.00	1/1 = 1.00	4/4 = 1.00	4/4 = 1.00	6/7 = 0.86	6/7 = 0.86
The Scholar's Bookshelf	1/1 = 1.00	1/1 = 1.00	1/1 = 1.00	3/3 = 1.00	3/3 = 1.00	14/18 = 0.78	14/17 = 0.83

Table 4. Experimental results for each site in *Music Shopping* domain.

Datasets / Metrics	Form-Domain Associations			Field-Attribute Associations		Text-Field Associations	
	Precision	Recall	Precision	Precision	Recall	Precision	Recall
			SubmittedForms				
Amazon.com Music	1/1 = 1.00	1/1 = 1.00	1/1 = 1.00	4/4 = 0.00	4/4 = 1.00	14/16 = 0.88	14/14 = 1.00
Barnes&Noble Music	1/1 = 1.00	1/1 = 1.00	1/1 = 1.00	4/4 = 1.00	4/4 = 1.00	7/10 = 0.70	7/7 = 1.00
Tower Records	1/1 = 1.00	1/1 = 1.00	1/1 = 1.00	6/6 = 1.00	6/6 = 1.00	19/19 = 1.00	19/19 = 1.00
Rough Trade	1/1 = 1.00	1/1 = 1.00	1/1 = 1.00	3/3 = 1.00	3/3 = 1.00	3/4 = 0.75	3/4 = 0.75
Sam Goody	1/1 = 1.00	1/1 = 1.00	1/1 = 1.00	3/3 = 1.00	3/3 = 1.00	8/10 = 0.80	8/8 = 1.00
Schott Musik International	1/1 = 1.00	1/1 = 1.00	1/1 = 1.00	3/3 = 1.00	3/3 = 1.00	12/14 = 0.86	12/12 = 1.00
A&B Sound	1/1 = 1.00	1/1 = 1.00	1/1 = 1.00	3/3 = 1.00	3/3 = 1.00	1/4 = 0.25	1/3 = 0.33
Collectors' Choice Music	1/1 = 1.00	1/1 = 1.00	1/1 = 1.00	5/5 = 1.00	5/5 = 1.00	5/8 = 0.63	5/5 = 1.00
CD Quest	1/1 = 1.00	1/1 = 1.00	1/1 = 1.00	2/2 = 1.00	2/2 = 1.00	10/11 = 0.91	10/10 = 1.00
AudibleFaith	1/1 = 1.00	1/1 = 1.00	1/1 = 1.00	4/4 = 1.00	4/4 = 1.00	13/14 = 0.93	13/13 = 1.00
Cyber Music Surplus	1/1 = 1.00	1/1 = 1.00	1/1 = 1.00	4/4 = 1.00	4/6 = 0.67	20/22 = 0.91	20/20 = 1.00
ClassicTrax.co.uk	1/1 = 1.00	1/1 = 1.00	1/1 = 1.00	4/4 = 1.00	4/4 = 1.00	8/9 = 0.88	8/9 = 0.89
MyMusic.com	1/1 = 1.00	1/1 = 1.00	1/1 = 1.00	4/4 = 1.00	4/6 = 0.67	9/11 = 0.82	9/9 = 1.00
Mojo Sounds	1/1 = 1.00	1/1 = 1.00	1/1 = 1.00	3/4 = 0.75	3/4 = 0.75	11/12 = 0.92	11/11 = 1.00
Looney Tunes	1/1 = 1.00	1/1 = 1.00	1/1 = 1.00	3/4 = 0.75	3/3 = 1.00	13/18 = 0.72	13/14 = 0.93
Buy Cd	1/1 = 1.00	1/1 = 1.00	1/1 = 1.00	3/3 = 1.00	3/3 = 1.00	9/14 = 0.64	9/9 = 1.00
Record Exchange	1/1 = 1.00	1/1 = 1.00	1/1 = 1.00	3/3 = 1.00	3/3 = 1.00	9/13 = 0.69	9/9 = 1.00
Musica Obscura	1/1 = 1.00	1/1 = 1.00	1/1 = 1.00	3/3 = 1.00	3/3 = 1.00	13/13 = 1.00	13/13 = 1.00
ProMusicFind.com	-	0/1 = 0.00	-	1/1 = 1.00	1/2 = 0.50	2/12 = 0.17	2/2 = 1.00
Ladyslipper Music	1/1 = 1.00	1/1 = 1.00	1/1 = 1.00	3/3 = 1.00	3/3 = 1.00	13/18 = 0.72	13/13 = 1.00

Table 5. Experimental results for each site in *Movies Shopping* domain.

Datasets / Metrics	Form-Domain Associations			Field-Attribute Associations		Text-Field Associations	
	Precision	Recall	Precision	Precision	Recall	Precision	Recall
			SubmittedForms				
Amazon.com DVD	2/2 = 1.00	2/2 = 1.00	2/2 = 1.00	10/10 = 1.00	10/10 = 1.00	39/41 = 0.95	39/39 = 1.00
The Spinning Image	1/1 = 1.00	1/1 = 1.00	1/1 = 1.00	4/4 = 1.00	4/4 = 1.00	13/14 = 0.93	13/13 = 1.00
Saregama	1/1 = 1.00	1/1 = 1.00	1/1 = 1.00	5/5 = 1.00	5/5 = 1.00	13/14 = 0.93	13/13 = 1.00
Docuseek	1/1 = 1.00	1/1 = 1.00	1/1 = 1.00	3/3 = 1.00	3/3 = 1.00	25/41 = 0.61	25/39 = 0.64
BollyVista.com	1/1 = 1.00	1/1 = 1.00	1/1 = 1.00	5/5 = 1.00	5/5 = 1.00	7/7 = 1.00	7/7 = 1.00
Movie Tickets.com	1/1 = 1.00	1/1 = 1.00	1/1 = 1.00	5/5 = 1.00	5/5 = 1.00	14/15 = 0.93	14/14 = 1.00
Yahoo! Movies	2/2 = 1.00	2/2 = 1.00	2/2 = 1.00	3/4 = 0.75	3/3 = 1.00	19/20 = 0.95	19/19 = 1.00
Half.com	1/1 = 1.00	1/1 = 1.00	1/1 = 1.00	4/4 = 1.00	4/4 = 1.00	5/8 = 0.63	5/5 = 1.00
GlobeAndMail.com	1/1 = 1.00	1/1 = 1.00	1/1 = 1.00	3/3 = 1.00	3/3 = 1.00	10/10 = 1.00	10/10 = 1.00
Kansas.com	1/1 = 1.00	1/1 = 1.00	1/1 = 1.00	3/3 = 1.00	3/3 = 1.00	9/9 = 1.00	9/9 = 1.00
Fandango	1/1 = 1.00	1/1 = 1.00	1/1 = 1.00	4/4 = 1.00	4/4 = 1.00	9/10 = 0.90	9/9 = 1.00
Sensasian.com	1/1 = 1.00	1/1 = 1.00	1/1 = 1.00	4/4 = 1.00	4/4 = 1.00	13/14 = 0.93	13/14 = 0.93
Blockbuster.com	1/1 = 1.00	1/1 = 1.00	1/1 = 1.00	4/4 = 1.00	4/4 = 1.00	16/18 = 0.89	16/17 = 0.94
Video*ezzy	1/1 = 1.00	1/1 = 1.00	1/1 = 1.00	3/3 = 1.00	3/3 = 1.00	9/9 = 1.00	9/9 = 1.00
Intelliflix	1/1 = 1.00	1/1 = 1.00	1/1 = 1.00	3/3 = 1.00	3/3 = 1.00	18/18 = 1.00	18/18 = 1.00
IGN.COM	-	0/1 = 0.00	-	2/2 = 1.00	2/2 = 1.00	65/65 = 1.00	65/70 = 0.93
Bestvideobuys	1/1 = 1.00	1/1 = 1.00	1/1 = 1.00	4/4 = 1.00	4/4 = 1.00	5/5 = 1.00	5/5 = 1.00
BrickFilms	1/1 = 1.00	1/1 = 1.00	1/1 = 1.00	3/3 = 1.00	3/4 = 0.75	16/22 = 0.73	16/16 = 1.00
Human Rights Film Dir.	1/1 = 1.00	1/1 = 1.00	1/1 = 1.00	2/2 = 1.00	2/3 = 0.67	0/9 = 0.00	0/9 = 0.00
Blockbuster UK	1/1 = 1.00	1/1 = 1.00	1/1 = 1.00	4/4 = 1.00	4/4 = 1.00	12/14 = 0.86	12/14 = 0.86

6 Related Work

In recent years, several works have addressed the problem of accessing the Hidden Web using a variety of approaches.

The system more similar to ours is HiWE [17]. HiWE is a task-specific crawler able to automatically recognizing and filling in forms relevant to a given data-collecting task. As well as DeepBot, HiWE also uses visual distance measures to find the texts associated to each field in a form, and textual similarity measures to match form fields and domain attributes.

When analyzing forms, HiWE only associates one text to each form field. The text is chosen in the following way: first, HiWE finds the four closest texts to the field; second, it chooses one of them according to a set of heuristics. These heuristics take into account the relative position of the candidate texts with respect to the field (texts at the left and at the top are privileged), and their font sizes and styles.

To learn how to fill in a form, HiWE matches the text associated with each form field and the labels associated to the attributes defined in its LVS table (a concept that plays a similar role to our domain definitions). In this process, HiWE has the following restriction: it requires the LVS table to contain an attribute definition matching with each unbounded form field.

Now we discuss the differences between HiWE and our system. The process followed by DeepBot has several advantages:

- DeepBot may use a form, even though it has some fields that do not match any attribute of the domain. For instance, the domain definition in Fig. 2 does not have any attribute matching with the “*Publisher*” field in Fig. 3a, but DeepBot would be able to use the form anyway.
- DeepBot correctly detects when a field has more than one associated text; this can result in better accuracy when matching form fields and domain attributes.
- In addition, the decision of assigning a text to a field is not based only on conditions “local” to the field: the context provided by the whole form is also taken into account in our heuristics. For instance, in our example form of Fig. 3a, HiWE would erroneously assign the text “*Hardcover*” to the second *radio button* element (f_{52}), since the text is the closest one and, in addition, it is located at the left of the field. Nevertheless, our system correctly assigns the text “*e-Books & Docs*” to f_{53} , “*Paperback*” to f_{52} and “*Hardcover*” to f_{51} . That is because DeepBot guarantees that every field will be associated to at least one text (the assignments made by HiWE would let f_{51} without any assigned text or would assign the same text to two fields).
- Finally, another important advantage of DeepBot is that it fully supports Javascript sources, while HiWE does not.

Reference [3] presents another system for domain-specific crawling of the Hidden Web. Nevertheless, they only deal with *full text* search forms; these forms have a single field allowing search by keyword on unstructured document collections. In turn, our system focuses on the multi-attribute forms typically used to query structured data.

Reference [15] addresses the problem of automatically generating keyword queries to crawl all the content behind a web form. New techniques are proposed to automatically generate new search keywords from previous results, and to prioritize them in order to retrieve the full content behind the form, using the minimum number of queries. The ability to automatically generate new queries from the results of previous ones would be an interesting new feature for DeepBot, so this work is complementary to ours. Nevertheless, the presented techniques would need to be significantly adapted since they focus only on keyword search forms and do not deal with multi-attribute forms.

The problem of extracting the full content behind a web form has been also addressed in [14]. Nevertheless, this system does not deal with forms requiring *textbox* fields to be filled in.

The Hidden Web can also be accessed using the *meta-search* paradigm instead of the *crawling* paradigm. In *meta-search* systems, a query from the user is automatically redirected to a set of underlying relevant sources, and the obtained results are integrated to return a unified response.

The meta-search approach is more lightweight than the crawling approach, since it does not require indexing the content from the sources; it also guarantees up to date data. Nevertheless, users will get higher time responses since the sources are queried

in real-time. Some systems using the meta-search approach are MetaQuerier [6], [9], [20], WISE-Integrator [10], [11] and QProber [8], [12], [13].

Finally, several techniques [2], [18], [19] have been proposed to automatically extract structured data from web pages conforming to a certain template. Since a substantial portion of the Hidden Web is composed of forms providing access to underlying structured databases, these techniques are also complementary to our work.

7 Conclusions

In this paper, we have described the architecture of *DeepBot*, a crawling system able to access the contents of the Hidden Web. We have focused on the techniques used to access the content behind web forms (server-side Deep Web). Our approach is based on a set of domain definitions, each one describing a data-collecting task. From the domain definition, the system uses several heuristics, based on visual distance and text similarity measures, to automatically identifying relevant query forms and learning how to execute queries on them. We have tested our techniques for several real-world data-collecting tasks, obtaining a high degree of effectiveness.

References

1. Álvarez, M., Pan, A., Raposo, J., Hidalgo, J. Crawling Web Pages with Support for Client-Side Dynamism. Published in Lecture Notes in Computer Science 4016, pp. 252-262, 2006. Issue corresponding to Proceedings of the 7th International Conference on Web Age Information Management. 2006.
2. Arasu, A. and Garcia-Molina, H. Extracting Structured Data from Web Pages. In Proceedings of the ACM SIGMOD International Conference on Management of data. 2003.
3. Bergholz, A., Chidlovskii, B. Crawling for Domain-Specific Hidden Web Resources. In Proceedings of the 4th Int. Conference on Web Information Systems Engineering. 2003.
4. Bergman, M. The Deep Web. Surfacing Hidden Value. <http://brightplanet.com/technology/deepweb.asp>. 2001.
5. C.-C. Chang, K., He, B., Patel, M., Zhang, Z. Structured Databases on the Web: Observations and Implications. SIGMOD Record, 33(3). 2004.
6. C.-C. Chang, K., He, B., Zhang, Z. MetaQuerier over the Deep Web: Shallow Integration Across Holistic Sources. In Proceedings of the VLDB Workshop on Information Integration on the Web. 2004.
7. Cohen, W., Ravikumar, P., Fienberg, S. A Comparison of String Distance Metrics for Name-Matching Tasks. In Proceedings of IJCAI-03 Workshop. 2003.
8. Gravano, L., Ipeirotis, P., Sahami, M. QProber: A System for Automatic Classification of Hidden-Web Databases. In ACM Transactions on Information Systems, vol. 21(1), 2003.
9. He, B., C.-C. Chang, K. Automatic Complex Schema Matching across Web Query Interfaces: A Correlation Mining Approach. In ACM Transactions on Database Systems (TODS), vol. 31(1), 2006

10. He, H., Meng, W., Yu, C., and Wu, Z. Automatic Integration of Web Search Interfaces with WISE-Integrator. In VLDB Journal, Vol.13, No.3, pp.256-273. 2004. (Special Issue for Best Papers of VLDB 2003)
11. He, H., Meng, W., Yu, C., and Wu, Z. WISE-Integrator: A System for Extracting and Integrating Complex Web Search Interfaces on the Deep Web. In Proceedings of the 31th VLDB Conference. 2005.
12. Ipeirotis P., Gravano L. Distributed Search over the Hidden Web: Hierarchical Database Sampling and Selection. In Proceedings of the 28th International Conference on Very Large Databases. 2002.
13. Ipeirotis, P., Ntoulas A., Cho, J., Gravano, L. Modeling and Managing Content Changes in Text Databases. In Proceedings of the 21st IEEE International Conference on Data Engineering, 2005
14. Liddle, S., Embley, D., Scott, Del., Yau Ho, Sai. Extracting Data Behind Web Forms. Proceedings of the 28th Intl. Conference on Very Large Databases. 2002.
15. Ntoulas, A., Zerkos, P., Cho, J. Downloading Textual Hidden Web Content Through Keyword Queries. In Proceedings of the 5th ACM/IEEE Joint Conference on Digital Libraries. 2005.
16. Pan, A., Raposo, J., Álvarez, M., Hidalgo, J. and Viña, A. Semi-Automatic Wrapper Generation for Commercial Web Sources. In Proceedings of IFIP WG8.1 Working Conference on Engineering Information Systems in the Internet Context. 2002.
17. Raghavan S., Garcia-Molina, H. Crawling the Hidden Web. Technical Report 2000-36, Computer Science Department, Stanford University, December 2000. Available at <http://dbpubs.stanford.edu/pub/2000-36>
18. Wang J., Lochovsky F. Data Extraction and Label Assignment for Web Databases. In Proceedings of the 12th World Wide Web Conference. 2003
19. Zhai, Y., Liu, B. Structured Data Extraction from the Web Based on Partial Tree Alignment. In IEEE Transactions on Knowledge and Data Engineering Journal, vol 18, n°12. 2006.
20. Zhang, Z., He, B., C.-C. Chang, K. Light-weight Domain-based Form Assistant: Querying Web Databases On the Fly. In Proceedings of the 31st Very Large Data Bases Conference, 2005.