

# A Task-specific Approach for Crawling the Deep Web

Manuel Álvarez, Juan Raposo, Fidel Cacheda and Alberto Pan

**Abstract**— There is a great amount of valuable information on the web that cannot be accessed by conventional crawler engines. This portion of the web is usually known as the Deep Web or the Hidden Web. Most probably, the information of highest value contained in the deep web, is that behind web forms. In this paper, we describe a prototype hidden-web crawler able to access such content. Our approach is based on providing the crawler with a set of *domain definitions*, each one describing a specific data-collecting task. The crawler uses these descriptions to identify relevant query forms and to learn to execute queries on them. We have tested our techniques for several real world tasks, obtaining a high degree of effectiveness.

**Index Terms**—Crawler, Deep Web, HTML Forms, Server-Side.

## I. INTRODUCTION

Crawlers are software programs that automatically traverse the web, retrieving pages to build a searchable index of their content. Conventional crawlers receive as input a set of "seed" pages and recursively obtain new ones by locating and traversing their outbound links.

Crawling techniques have led the construction of highly successful commercial web search engines. Nevertheless, conventional web crawlers cannot access to a significant fraction of the web, which is usually called the "hidden web" or the "deep web". 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 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.

Several works have tried to characterize the hidden web [4],

[5]. Among their findings, we can point out that the "hidden web" is substantially larger than the publicly indexable web and that hidden web pages (specially the ones accessed through query forms) usually contain data of higher quality and with a higher degree of structure.

In addition, the hidden web size is growing rapidly: in the four years lapse between the two cited studies, the estimated size of the hidden web has been increased between 3 and 7 times [5].

To address the problem of crawling the hidden web, we have built a prototype system called *DeepBot*. It has the following features:

- *DeepBot's* crawling processes are based on automated "mini web browsers", built by using standard browser APIs (our current implementation is based on the Microsoft Internet Explorer browser). This enables our system to deal with client-side scripting code, session mechanisms, managing redirections, and other complexities related with the client-side hidden web.
- 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.

This paper briefly overviews the architecture of *DeepBot* and describes in detail the techniques it uses for accessing the server-side hidden web. The techniques used to deal with the client-side deep web were described in greater detail in [1].

The rest of the paper is organized as follows. Section II overviews the architecture of *DeepBot* and the main components that participate in accessing the server-side hidden web. Section III describes the domain definitions used to specify a data collection task. Section IV is the core of the paper; it describes how *DeepBot* detects query forms relevant to a certain task and how it learns to execute queries on them. Section V describes our experiments with the system. Finally, section VI discusses related work and section VII concludes the paper.

## II. OVERVIEW / ARCHITECTURE

As well as in conventional crawling engines, 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

Manuscript received April 28, 2006. This research was supported in part by the Spanish Ministry of Education and Science under project TSI2005-07730. Alberto Pan's work was supported in part by the "Ramón y Cajal" programme of the Spanish Ministry of Education and Science.

M. Álvarez, J. Raposo, F. Cacheda and A. Pan are with the Department of Information and Communications Technology, University of A Coruña, Campus de Elviña s/n. 15071 A Coruña, Spain (e-mail: {mad, jrs, fidel, apan}@udc.es).

machines.

The crawler is initialized with a list of routes. Each crawling process picks a route from the list, downloads its associated document and analyzes it in order to obtain new routes from its anchors, which are then added to the master list. The process ends when there are no routes left or when a specified depth level is reached. 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. This allows a crawling process to access a URL added by other crawling process (even if the original process was running in another machine).
- 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 browser technologies) 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 instance, to access to the page behind an anchor, a conventional crawler would obtain the URL from the *href* attribute, using it to issue an HTTP request. Instead, our system simply generates a *click* browser-event on the anchor. This allows the crawling processes to forget about complexities such as client-side scripting (e.g. Javascript) or complex redirections, in the same way a human user of a web browser is not bothered about those issues. For specifying a navigation sequence in the automated mini-browsers, we have created NSEQL [15], 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. For instance, NSEQL includes commands for actions such as generating 'click' events on any element of a page (*anchors, images, buttons, ...*), filling in HTML forms, etc.

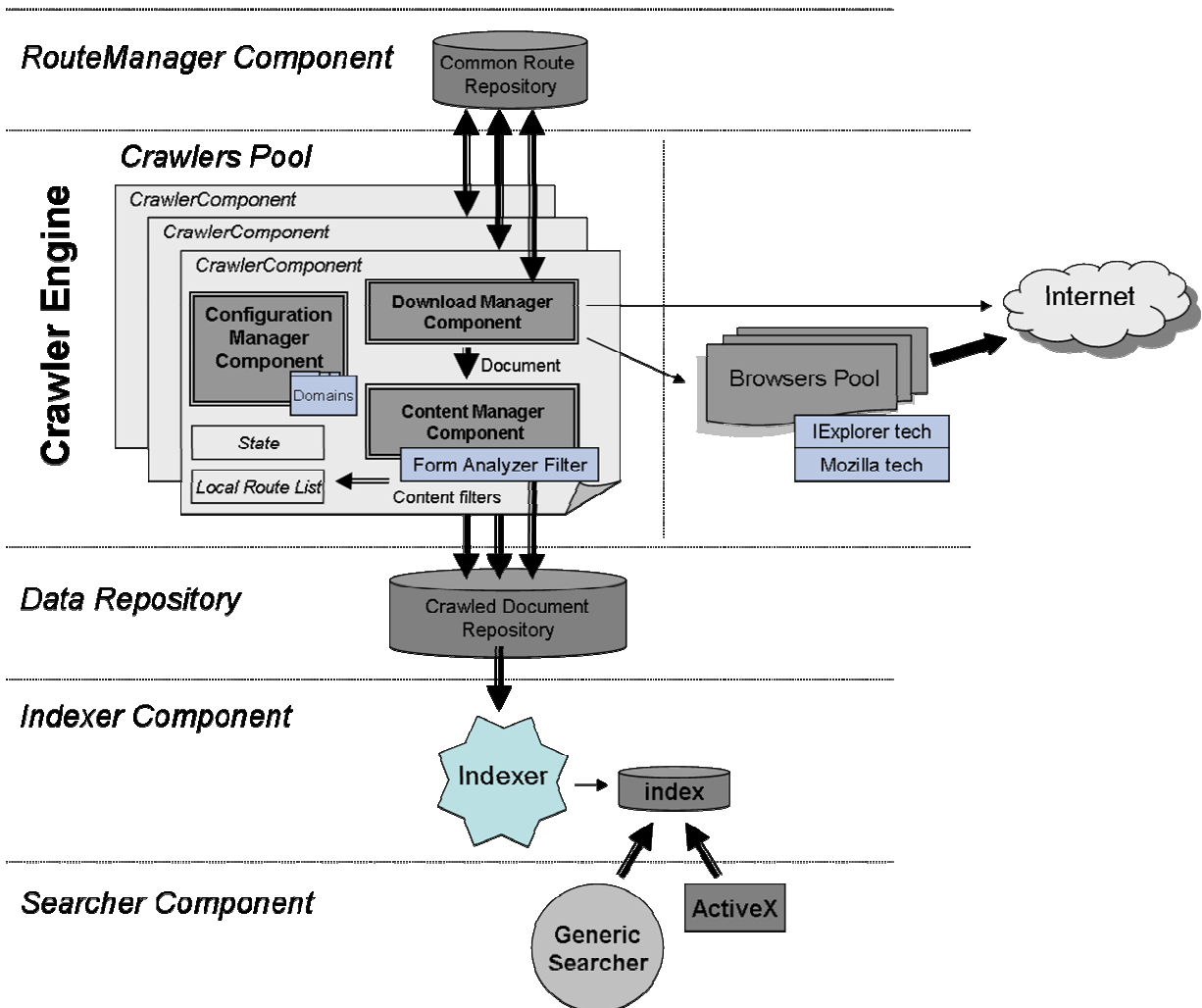


Fig. 1. Crawler Architecture

- 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.

#### A. Crawler Architecture

The architecture of the system is shown in Fig. 1. When the crawler engine starts, it reads its configuration parameters from the *Configuration Manager* component.

These parameters include the ones typically needed by a conventional crawler engine: the list of initial routes to begin the crawl, the desired navigation depth for each initial route, download handlers for different kinds of documents, content filters, a list of regular expressions representing URLs to be included and excluded from the crawling, etc.

In addition, the crawler is also configured with a set of *domain definitions*. As we will see in section III, each domain definition describes a specific data-collection task in the server-side deep web.

The *Route Manager* is responsible for maintaining the master list of *routes* to access; all crawlers share this list. Once the crawling processes start, each one picks a route from the Route Manager. It is important to notice that each crawling process can be executed either locally or remotely to the server, thus allowing for distributed crawling. As we have already remarked, each crawling process is a mini web-browser able to execute NSEQL sequences.

Then the crawling process loads the session object associated to the route and downloads the associated document (it uses the *Download Manager Component* to choose the right handler for it, such as HTML, PDF, MS Word, etc).

The content from each downloaded document is then analyzed by using the *Content Manager Component*. This component specifies two chains of filters. The first chain is used to decide if the document is considered relevant and, therefore, if it should be stored and/or indexed. For instance, the system includes filters, which allow checking if the document verifies a keyword-based boolean query in order to decide whether to store/index it or not.

The second chain of filters is used to post-process the document. For instance, this chain includes filters to extract the useful content from HTML pages and to generate a short document summary. Nevertheless, the main function of this second chain is obtaining new routes from the analyzed documents and adding them to the master list. There are two filters involved in this task:

- The *Obtain Links* filter selects all the anchors in the page and generates a new route for each one. Since scripting languages can dynamically generate and remove anchors in response to user actions (e.g. pop-up menus), this involves some complexities for a hidden-web crawler (see [1]). It is also possible to specify a certain

regular expression that the URL of the anchors must verify in order to be considered.

- The *Form Analyzer* filter analyzes each form in the page and determines if it is relevant for any of the pre-configured domain definitions. In the case a form is considered relevant, a new route will be added for each query specified by the domain definition. Subsequent sections will provide more detail about this filter.

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<sup>1</sup>). The crawler generates an XML file for each crawled document, including meta-information such as its URL and the NSEQL sequence needed to access it.

The NSEQL sequence is used by another component of the system architecture: *the ActiveX for automatic navigation Component*. This component receives as a parameter a NSEQL program, downloads itself into the user browser and makes it execute the given navigation sequence. In our system, this is used to solve the problem of accessing to documents at a later time. When the user executes a search against the index, the list of answers may contain some results that cannot be directly accessed by using its URL, due to session issues. In that case, the anchors associated to those results in the response listing will invoke the ActiveX component passing as parameter the NSEQL sequence associated to the page. Then, if the users click on the anchor, the ActiveX will make their browser automatically navigate to the desired page.

### III. DOMAIN DEFINITIONS

In this section, we describe the domain definitions used to describe 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*  $s_i$ .
- 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  $\mu$ .

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. In addition, they also detected a

<sup>1</sup> <http://lucene.apache.org>

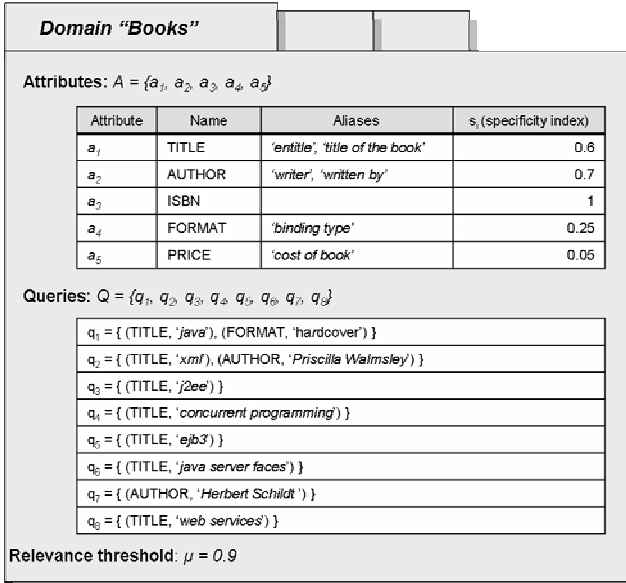


Fig. 2. Example domain definition for Books

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, we could have the following specificity indexes:

- 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 allowing queries for the PRICE attribute could be a query form allowing to search any kind of product, not only books.

Finally, the domain also includes a *relevance threshold*  $\mu$ . 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.

#### IV. PROCESSING FORMS WITH THE FORM ANALYZER

In this section, we describe how the crawler processes each found form (see Fig. 3). 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 (described in subsection A).
- By using the output of the previous step, the system determines if the form is relevant with respect to the domain (see subsection B).
- 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 (see subsection C).

The following sub-sections detail each of these steps.

##### A. 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. The method we use consists of the following 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<sup>2</sup> 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

<sup>2</sup> We consider the radio button and checkbox elements with the same value for the *name* attribute as a single field. The system tries to associate texts for each *radio button* or checkbox, and texts for the compound field they represent (group of *radio buttons* or *checkboxes*).

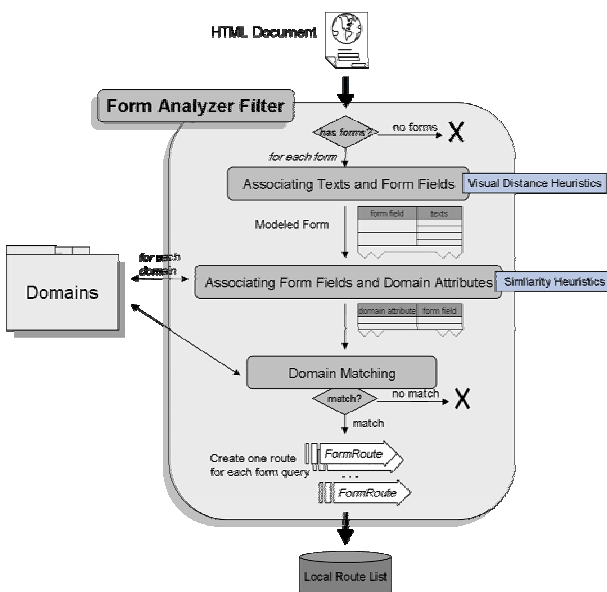


Fig. 3. Form Analyzer Architecture

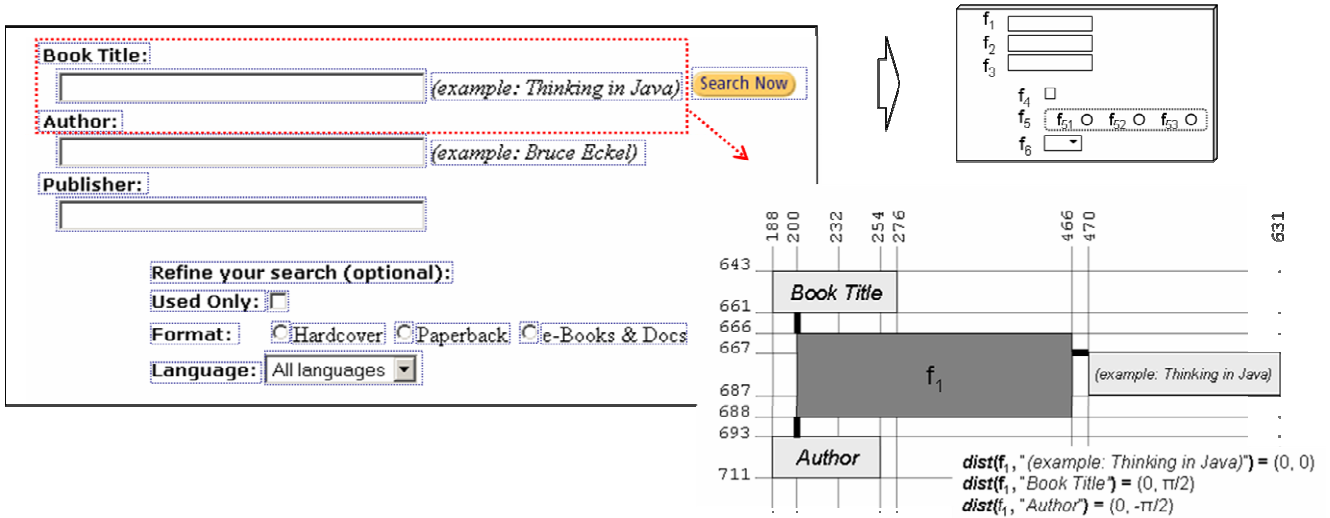


Fig. 4a. Example query form and visual distances and angles for field  $f_1$

rectangle enclosing  $f$  and a rectangle enclosing  $t$ . If  $t$  is into an HTML table cell, and it is the unique text inside that table cell, then the coordinates of the table cell rectangle are assigned to  $t$ .

2. We obtain the minimum distance between both rectangles. This involves finding the shortest line joining any point in one rectangle and any point in the other; the distance will be the length of this line. 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. 4a 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. 4a 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. First, we add all the texts having the shortest distance  $d$  with respect to  $f$  to the list.
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 than the closest ones.
3. Texts with the same distance are ordered according to its angle (note that since our distances are rather measured in coarse-grained units than pixels, it is relatively usual to have several texts at the same distance from the same field). 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.

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. Note that although there may be texts in the form which are semantically related to more than one field (e.g. the text “Refine your search:” in Fig. 4a), those semantic associations will typically be irrelevant for our purposes; because these texts are related to more than one field, they usually do not describe precisely any of them.
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. 4b shows the process for the example form of Fig. 4a. For each field<sup>3</sup> 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. For instance, for the field  $f_1$  the final associated texts are “(example: Thinking in Java)” and “Book

<sup>3</sup> Note how the system models the *Format* ‘checkbox’ field as a field with three subfields.  $f_3$  refers to the whole set of *checkboxes* while  $f_{31}$ ,  $f_{32}$  and  $f_{33}$  refer to individual *checkboxes*.

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

Fields	Texts	(dist, $\theta$ )
$f_{51}$	√ Hardcover	(0, 0)
	Used Only:	(0, $3\pi/4$ )
	Language:	(0, $-3\pi/4$ )
	Format:	(1, $\pi$ )
	Refine your search (optional):	(1, $\pi/2$ )
$f_{52}$	√ Hardcover	(0, $\pi$ )
	Paperback	(0, 0)
$f_{53}$	√ Paperback	(0, $\pi$ )
	e-Books & Docs	(0, 0)
$f_6$	√ Language:	(0, $\pi$ )
	Hardcover	(0, $\pi/2$ )
	Paperback	(0, $\pi/2$ )
	Format:	(1, $\pi$ )
	Used Only:	(1, $\pi/2$ )
	Refine your search (optional):	(3, $\pi/2$ )

$f_1$  [ (example: Thinking in Java) ] [ Book Title: ]  
 $f_2$  [ (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. 4b. Texts associated to each field in the form of Fig. 4a

Title:”.

*Associating form fields and domain attributes.* At this step we try to detect the form fields which correspond to attributes of the target domain.

At this stage, 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*,
- *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 page<sup>4</sup> 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

<sup>4</sup> 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.

Assignment	Form Field	Domain Attribute	$c_i$ (confidence)
$A_1$	$f_1$	$a_1 = \text{TITLE}$	0.71
$A_2$	$f_2$	$a_2 = \text{AUTHOR}$	1
	$f_3$	(unassigned)	
	$f_4$	(unassigned)	
$A_3$	$f_5$	$a_4 = \text{FORMAT}$	1
	$f_6$	(unassigned)	



**Assignments = { $A_1, A_2, A_3$ }**

TITLE

AUTHOR

$f_3$

$f_4$

FORMAT   $f_{51}$    $f_{52}$    $f_{53}$

$f_6$

Fig. 5. Assignments obtained for the form in Fig.4a, using the domain definition shown in Fig.2

with  $f$  in the page using the methods seen in the previous section; we notate them as  $\{f\_text_1, \dots, f\_text_m\}$ .

- 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  $a$   $\{a\_value_1, \dots, a\_value_q\}$ .
- 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.

$$sim_1 = \max \left\{ \text{textSim}(a\_alias_i, f\_text_j)_{i=1..n, j=1..m} \right\} \quad (1)$$

- 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  $a\_value_r, r=1..q$ .

$$sim_2 = \sum_{r=1..q} \left( \max \left\{ \text{textSim}(a\_value_r, f\_value_k)_{k=1..p} \right\} / q \right) \quad (2)$$

- 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. For instance, if the field  $f_2$  were the best option for two attributes  $a_1$  and  $a_2$ , and we also had  $f_1$  as an additional option for  $a_1$ , we would make the assignments  $f_1-a_1$  and  $f_2-a_2$  instead of the sole assignment  $f_2-a_2$ .

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 participating in the

assignment.

Fig. 5 shows the assignments obtained for the form in Fig. 4a, using the domain definition shown in Fig. 2.

### B. 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*  $\mu$ . That is, the system checks if the inequality in (3) is verified.

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

For instance, considering the domain definition shown in Fig. 2, and the assignments shown in Fig. 5, we would obtain (4).

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

### C. 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.
- 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:

"Books Shopping"		
Music		
Attributes		
Attribute Name	Aliases	s <sub>i</sub> (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: $\mu = 0.9$		

Books		
"Music Shopping"		
Attributes		
Attribute Name	Aliases	s <sub>i</sub> (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: $\mu = 0.9$		

Fig. 6. Domain definitions: Books and Music

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

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.
3. The anchors obtained in the previous step 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', ...).
4. 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.
5. If all the previous steps fail, the system generates a SUBMIT event on the form.

## V. EXPERIENCE

### A. Experimental Setting

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

Table 1. Basic and advanced datasets for Books Shopping and Music Shopping domains

Domain	Books Shopping	Music Shopping
Basic	Amazon.com - <a href="http://www.amazon.com">http://www.amazon.com</a>	Amazon.com Music - <a href="http://www.amazon.com">http://www.amazon.com</a>
	Barnes&Noble - <a href="http://www.barnesandnoble.com">http://www.barnesandnoble.com</a>	Barnes&Noble Music - <a href="http://www.barnesandnoble.com">http://www.barnesandnoble.com</a>
	Powell's Books - <a href="http://www.powells.com">http://www.powells.com</a>	Tower Records - <a href="http://www.towerrecords.com">http://www.towerrecords.com</a>
	eCampus.com - <a href="http://www.ecampus.com">http://www.ecampus.com</a>	Rough Trade - <a href="http://www.roughtrade.com">http://www.roughtrade.com</a>
	Dymocks Booksellers - <a href="http://www.dymocks.com.au">http://www.dymocks.com.au</a>	Sam Goody - <a href="http://www.samgoody.com">http://www.samgoody.com</a>
	Tattered Cover Bookstore - <a href="http://www.tatteredcover.com">http://www.tatteredcover.com</a>	Schott Musik International - <a href="http://www.schott-music.com">http://www.schott-music.com</a>
	BookFinder4U.com - <a href="http://www.bookfinder4u.com">http://www.bookfinder4u.com</a>	A&B Sound - <a href="http://catalog2.absound.ca">http://catalog2.absound.ca</a>
	Cody's Books - <a href="http://www.codysbooks.com">http://www.codysbooks.com</a>	Collectors' Choice Music - <a href="http://www.ccmusic.com">http://www.ccmusic.com</a>
	Daedalus Books&Music - <a href="http://www.daedalusbooks.com">http://www.daedalusbooks.com</a>	CD Quest - <a href="http://www.cdquest.com">http://www.cdquest.com</a>
	Bolen - <a href="http://www.bolen.bc.ca">http://www.bolen.bc.ca</a>	AudibleFaith - <a href="http://www.audiblefaith.com">http://www.audiblefaith.com</a>
	Advanced	The Book Pl@ce - <a href="http://www.thebookplace.com">http://www.thebookplace.com</a>
Blackwell's Bookshop - <a href="http://bookshop.blackwell.co.uk">http://bookshop.blackwell.co.uk</a>		ClassicTrax.co.uk - <a href="http://www.classictrax.ltd.uk">http://www.classictrax.ltd.uk</a>
The American Book Center - <a href="http://www.abc.nl">http://www.abc.nl</a>		MyMusic.com - <a href="http://www.mymusic.com">http://www.mymusic.com</a>
Oxbow Books - <a href="http://www.oxbowbooks.com">http://www.oxbowbooks.com</a>		Mojo Sounds - <a href="http://www.mojosounds.com">http://www.mojosounds.com</a>
Strand Book Store - <a href="http://www.strandbooks.com">http://www.strandbooks.com</a>		Looney Tunes - <a href="http://www.looneytunescds.com">http://www.looneytunescds.com</a>
Globe Pequot Press - <a href="http://www.globepequot.com">http://www.globepequot.com</a>		Buy Cd - <a href="http://www.buycd.com">http://www.buycd.com</a>
Northshire Bookstore - <a href="http://www.northshire.com">http://www.northshire.com</a>		Record Exchange - <a href="http://www.buymusichere.net">http://www.buymusichere.net</a>
Green Apple Books&Music - <a href="http://www.greenapplebooks.com">http://www.greenapplebooks.com</a>		Musica Obscura - <a href="http://www.musicaobscura.com">http://www.musicaobscura.com</a>
Thomson Gale - <a href="http://www.gale.com">http://www.gale.com</a>		ProMusicFind.com - <a href="http://www.promusicfind.com">http://www.promusicfind.com</a>
The Scholar's Bookshelf - <a href="http://www.scholarsbookshelf.com">http://www.scholarsbookshelf.com</a>		Ladyslipper Music - <a href="http://www.ladyslipper.org">http://www.ladyslipper.org</a>



The process for creating the domain definitions was the following: for each domain, we manually explored 10 sites at random, from the respective Yahoo Directory<sup>5</sup> 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. 6.

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 *Basic*, 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.

### B. Metrics

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.

We defined the following metrics:

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

$$Precision_{TextFieldA} := \frac{|TextFieldA_{DeepBot} \cap TextFieldA_{Real}|}{|TextFieldA_{DeepBot}|} \quad (5)$$

$$Recall_{TextFieldA} := \frac{|TextFieldA_{DeepBot} \cap TextFieldA_{Real}|}{|TextFieldA_{Real}|}$$

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

$$Precision_{FieldAttributeA} := \frac{|FieldAttributeA_{DeepBot} \cap FieldAttributeA_{Real}|}{|FieldAttributeA_{DeepBot}|} \quad (6)$$

$$Recall_{FieldAttributeA} := \frac{|FieldAttributeA_{DeepBot} \cap FieldAttributeA_{Real}|}{|FieldAttributeA_{Real}|}$$

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

$$Precision_{FormDomainA} := \frac{|FormDomainA_{DeepBot} \cap FormDomainA_{Real}|}{|FormDomainA_{DeepBot}|} \quad (7)$$

$$Recall_{FormDomainA} := \frac{|FormDomainA_{DeepBot} \cap FormDomainA_{Real}|}{|FormDomainA_{Real}|}$$

$$Precision_{SubmittedForms} := \frac{|SubmittedForms_{DeepBot}|}{|FormDomainA_{DeepBot} \cap FormDomainA_{Real}|}$$

### C. Experimental Results

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

Table 2. Experimental results

Metrics / Datasets	Books Shopping			Music Shopping		
	Basic	Advanced	Global	Basic	Advanced	Global
<b>Form-Domain Associations</b>						
$Precision_{FormDomainA}$	13/13 = 1.00	11/11 = 1.00	24/24 = 1.00	10/10 = 1.00	9/9 = 1.00	19/19 = 1.00
$Recall_{FormDomainA}$	13/13 = 1.00	11/11 = 1.00	24/24 = 1.00	10/10 = 1.00	9/10 = 0.90	19/20 = 0.95
$Precision_{SubmittedForms}$	13/13 = 1.00	11/11 = 1.00	24/24 = 1.00	10/10 = 1.00	9/9 = 1.00	19/19 = 1.00
<b>Field-Attribute Associations</b>						
$Precision_{FieldAttributeA}$	54/55 = 0.98	50/50 = 1.00	104/105 = 0.99	37/37 = 1.00	31/33 = 0.94	68/70 = 0.97
$Recall_{FieldAttributeA}$	54/54 = 1.00	50/53 = 0.94	104/107 = 0.97	37/37 = 1.00	31/37 = 0.84	68/74 = 0.92
<b>Text-Field Associations</b>						
$Precision_{TextFieldA}$	129/142 = 0.91	101/137 = 0.73	230/279 = 0.82	93/110 = 0.83	107/132 = 0.81	199/242 = 0.82
$Recall_{TextFieldA}$	129/132 = 0.98	101/127 = 0.79	230/259 = 0.88	92/94 = 0.98	107/109 = 0.98	199/203 = 0.98

<sup>5</sup> <http://dir.yahoo.com>

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 SubmittedForms	Precision	Recall	Precision	Recall
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 SubmittedForms	Precision	Recall	Precision	Recall
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

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 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.

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.

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 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 and 4 show the obtained experimental results for each website.

## VI. 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 [14]. 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. 4a, 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. 4a, 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 [13] 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 [12]. 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], [18] and QProber [8], [10], [11].

Finally, several techniques [2], [16], [17] 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

## VII. 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.

## VIII. REFERENCES

- [1] Alvarez, M., Pan, A., Raposo, J., Hidalgo, J. Crawling Web Pages with Support for Client-Side Dynamism. To be published in Lecture Notes in Computer Science 4016, pp. 252-262, 2006. Issue corresponding to Proceedings of the 7<sup>th</sup> International Conference on Web Age Information Management (WAIM06). 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 International Conference on Web Information Systems Engineering (WISE). 2003.
- [4] Bergman, M. The Deep Web. Surfacing Hidden Value. <http://brightplanet.com/technology/deepweb.asp>. 2001.
- [5] Chen-Chuan Chang, K., He, B., Patel, M., Zhang, Z. Structured Databases on the Web: Observations and Implications. SIGMOD Record, 33(3). September 2004.
- [6] Chen-Chuan 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 (VLDB-IIWeb'04). 2004.
- [7] Cohen, W., Ravikumar, P., Fienberg, S. A Comparison of String Distance Metrics for Name-Matching Tasks. In Proceedings of IJCAI-03 Workshop (IIWeb-03). 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), Jan. 2003
- [9] He, B., Chen-Chuan Chang, K. Automatic Complex Schema Matching across Web Query Interfaces: A Correlation Mining Approach. In ACM Transactions on Database Systems (TODS), vol. 31(1), March 2006
- [10] 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 (VLDB2002). 2002.
- [11] 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 (ICDE 2005), 2005
- [12] Liddle, S., Embley, D., Scott, Del., Yau Ho, Sai. Extracting Data Behind Web Forms. Proceedings of the 28th Intl. Conference on Very Large Databases (VLDB2002). 2002.
- [13] 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 (JCDL05). 2005.
- [14] 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>
- [15] 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 (EISIC). 2002.
- [16] Wang J., Lochovsky F. Data Extraction and Label Assignment for Web Databases. In Proceedings of the 12th World Wide Web Conference. 2003
- [17] Zhai, Y., Liu, B. Web Data Extraction Based on Partial Tree Alignment. In Proceedings of the 2005 World Wide Web Conference (WWW2005). ACM Press. 2005
- [18] Zhang, Z., He, B., Chen-Chuan Chang, K. Light-weight Domain-based Form Assistant: Querying Web Databases On the Fly. In Proceedings of the 31st Very Large Data Bases Conference (VLDB 2005), 2005.