

SG-WRAP: A Schema-Guided Wrapper Generator

Xiaofeng Meng *
Information School
Renmin University of China
Beijing, China
xfmeng@public.bta.net.cn

Hongjun Lu †
Hong Kong University
of Science & Technology
Hong Kong, China
luhj@cs.ust.hk

Haiyan Wang, Mingzhe Gu
Information School
Renmin University of China
Beijing, China

1. Introduction

With the development of the Internet, the World-Wide-Web has become everyone's invaluable information source. However, most of data on the Web is currently in the form of HTML pages, which is neither well-structured nor associated with schema. It is almost impossible to use such data efficiently. Web wrapper technology has been developed to transform unstructured /semi-structured data to semi-structured/structured data, which can be queried and analyzed using matured techniques developed in database and other fields. The main issues of wrapper generation include (1) to identify semantics of the data contained in an HTML document and, (2) to establish the mappings between its structure and its semantics. Various wrapper generation tools developed addressed these issues in different ways [1, 2, 3, 5, 4].

In this paper, we present a wrapper generator, SG-WRAP. It was design and developed based on the following observations. First, while user interaction is probably the best way to help the generation of wrappers for a specific HTML source quickly and accurately because of the diversity of HTML pages and limited semantic information expressed in HTML tags, user's efforts should be minimized as much as possible. Second, the ultimate goal for wrapper generation is to transform the original data into some structured one that is easy to consume, rather than to understand the structure of the original data. When a user gathers data from the Web, s/he must have her/his needs in her/his mind. It is often not necessary to generate wrappers for the entire HTML document. Therefore, SG-WRAP, adopts a novel, schema guided, approach for wrapper generation. With this approach, a user defines the schema of

data to be extracted from an HTML page in terms of data type descriptors (DTD) of XML [6]. The user also provides example mappings by associating data in the HTML page and elements in DTD. The system then induces the mapping rules and generates a wrapper that extracts data from the HTML page and produce an XML document conformed to the specified DTD.

Compared with other available wrapper generators, SG-WRAP approach has the following unique features. First, with the guidance of user-defined schema, the wrapper generated could be more accurate and better reflect the users requirements. For the same source page, different users can obtain different data based on their interests. Second, the extracted data follows the user-defined schema and is ready to populate to database directly. As such, the wrapper can be easily integrated into the data integration process. Third, minimum user interaction is required. the user only needs a few clicks to establish the correspondences between data elements in the HTML document and elements in the schema. The system will induce the data extraction rules from the given sample matching instances.

2. SG-WRAP: The System

SG-WRAP consists of five major components: Preprocessor, Schema Acquirer, Rule Generator, Rule Refiner and Wrapper Generator. The main user interface is shown in Figure 1. Module Preprocessor is responsible for setting up the environment for the system. It fetches the Web page using the URL (Figure 1, A) given by the user. The fetched HTML page is displayed (Figure 1, B). Another input of the system is a user-defined schema for extracted data, which is obtained by the Schema Acquirer and displayed on the screen (Figure 1, C). A user entered schema can be saved in a Schema Base to be used later, or shared for other sources from which the same kinds of data to be extracted. The fetched HTML page is parsed. Syntax errors, such as missing tags in the original HTML document are fixed during the parsing process. Internally

*This author's work is partially supported by a grant from the Natural Science Foundation of China (No. 60073014).

†This author's work is partially supported by the Research Grant Council of the Hong Kong Special Administrative Region, China (grant HKUST6092/99E) and a grant from the National 973 project of China (No. G1998030414).

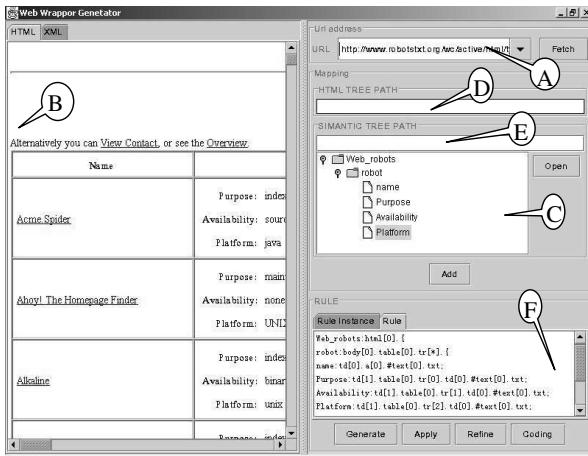


Figure 1. SG-WRAP: The main interface

the HTML page is represented as a tree using the document object model (DOM) [7], where a data item is a leaf node in the tree and its position in the document can be described by the path from the root to the leaf. As such, when the user highlights a string in the page shown, its path can be identified, which is displayed as an HTML tree path (Figure 1, D). Similarly, when the user clicks an element in the DTD, it is displayed as a semantic tree path in the schema (Figure 1, E). The Mapping Acquirer of the system captures those user clicks that associate strings in the HTML page and their corresponding elements in the DTD. For example, with the sample page in Figure 1, the following mapping can be captured

```
html.body.table[0].tbody[0].tr[0].a[0]
<=> Web_robots.robot.name "Alkaline"
```

where Alkaline is a string in the HTML page, `html.body.table[0].tbody[0].tr[0].a[0]` is the tree path in the HTML tree, and `Web_robots.robot.name` is the associated element in the DTD, respectively. In other words, from user interactions, the system obtains a set of instances of rules that map strings in the HTML page to elements in the schema. With this set of sample instances, the Rule Generator generates data extraction rule by an induction algorithm, which takes the list of mapping rules instances L as input and returns a candidate rule by incorporating the similar mapping rule instances into a new extraction rule. A sample data extraction rule for the HTML page shown in Figure 1 is as follows:

```
Web_robots:html{
  Robot: body.table[0].tbody[0].tr[] {
    Name: a[0].txt ;
    Purpose: table[0].tbody[0].tr[0].td[0].txt;
    Availability: table[0].tbody[0].tr[1].td[0].txt;
    Platform:table[0].tbody[0].tr[2].td[0].txt
  }
}
```

The data extraction rules are induced from a limited set of instances. In order to guarantee that the data extraction rule is applicable for the entire HTML document, SG-WRAP includes a refining process. The Rule Refiner generates an XML document by applying the induced rule on the input page. The generated document is displayed in Figure 1, window A for user to preview. From the displayed document, user determines whether the current extraction rule correctly extracts required data from the source page. If not, s/he can identify more mapping instances. The induction process will repeat and the Refiner will merge the previously generated rules with the refined rules. The system will display a new version of result to the user for checking. This refining process continues until the user is satisfied with the data extracted. The wrapper is generated based on the final data extraction rule by the Wrapper Generator.

3. Evaluations

Although wrapper generation work has been reported in the literature, there seem no standard ways to evaluate the performance of such systems. We conducted a series of experiments to evaluate the usability, correctness and efficiency of SG-WRAP. The usability tests selected a number of users to use the system. The results indicated that, with minimal introduction of the system, DTD definition and structure of HTML pages, even naive users could quickly generate wrappers without much difficulty. For correctness, we adapted the precision and recall metrics in information retrieval to data extraction. The results show that, with the refining process, the system can generate wrappers with very high accuracy. Finally, the efficiency tests indicated that the wrapper generation process is fast enough even with large size Web pages.

References

- [1] J. Hammer and *et al.* Template-based wrappers in the tsmis system. In *ACM SIGMOD*, pages 532–535, Tucson, Arizona, May 1997.
- [2] C. Knoblock and *et al.* Modeling web sources for information integration. In *AAAI*, pages 211–218, Madison, WI, 1998.
- [3] N. Kushmerick, D. Weil, and R. Doorenbos. Wrapper induction for information extraction. In *IJCAI*, pages 729–735, Nagoya, Japan, 1997.
- [4] L.Liu, C. Pu, and W. Han. XWRAP: An XML-enabled wrapper construction system for web information sources. In *ICDE*, pages 611–621, San Diego, CA, 2000.
- [5] A. Sahuguet and F. Azavant. WysiWyg web wrapper factory (W4F). In *WWW*, Toronto, Oct. 1999.
- [6] World Wide Web Consortium (W3C). Extensible markup language (XML1.0).
- [7] World Wide Web Consortium (W3C). The document object model. <http://www.w3.org/DOM>, 1998.