

DEPARTAMENTO DE SISTEMAS INFORMÁTICOS Y COMPUTACIÓN
UNIVERSIDAD POLITÉCNICA DE VALENCIA

P.O. Box: 22012 E-46071 Valencia (SPAIN)



Informe Técnico / Technical Report

Ref. No.: DSIC-II/08/07	Pages: 23
Title: The Web Verification Service WebVerdi-M	
Author(s): M. Alpuente, D. Ballis, M. Falaschi, P. Ojeda and D. Romero	
Date: 26 March 2007	
Keywords: web verification, maude, web service	

Vº Bº
Leader of research Group

Author(s)

The Web Verification Service WebVerdi-M ^{*}

M. Alpuente¹, D. Ballis², M. Falaschi³, P. Ojeda¹, and D. Romero¹

¹ DSIC, Universidad Politécnic de Valencia
Camino de Vera s/n, Apdo. 22012, 46071 Valencia, Spain.
{alpuente, pojeda, dromero}@dsic.upv.es

² Dip. Matematica e Informatica
Via delle Scienze 206, 33100 Udine, Italy.
demis@dimi.uniud.it

³ Dip. di Scienze Matematiche e Informatiche
Pian dei Mantellini 44, 53100 Siena, Italy.
moreno.falaschi@unisi.it

Abstract. In this work, we present the rewriting-based, Web verification service **WebVerdi-M**, which is able to recognize forbidden/incorrect patterns and incomplete/missing Web pages. **WebVerdi-M** relies on a powerful Web verification engine that is written in Maude, which automatically derives the error symptoms. Thanks to the AC pattern matching supported by Maude and its metalevel facilities, **WebVerdi-M** enjoys much better performance and usability than a previous implementation of the verification framework. By using the XML Benchmarking tool `xmlgen`, we develop some scalable experiments which demonstrate the usefulness of our approach.

1 Introduction

Web-site management is an arduous task today. While it is not so difficult to build a Web site, this task is still somehow a form of art, resulting in an increasing volume of information contained in ever-larger, complex Web sites, which is very difficult to keep up-to-date and correct. Recently, some sophisticated Web-site management tools have been proposed. These provide helpful facilities (see [4,5]), but unfortunately, they are mostly oriented to Web-site syntactic checking/restructuring so that they can do little by themselves to relieve the problem.

The automated management of data-intensive Web sites is an area to which rule-based technology has a significant potential to contribute. Web sites typically contain and integrate several bodies of data that are linked into a rich navigational structure. It is widely accepted today that declarative representations are the best way to specify the structural aspects of Web sites as well as many forms of Web-site content. As an additional advantage, rule-based languages such

^{*} This work has been partially supported by the EU (FEDER) and Spanish MEC TIN-2004-7943-C04-02 project, the Generalitat Valenciana under grant GV06/285, and Integrated Action Hispano-Alemana HA2006-0007. Daniel Romero is also supported by ALFA grant LERNet AML/19.0902/97/0666/II-0472-FA.

as Maude [12] offer an extremely powerful, rewriting-based “reasoning engine” where the system transitions are represented/derived by rewrite rules indicating how a configuration is *transformed* into another.

In previous work [3,6], we proposed a rewriting-based approach to Web-site verification and repair. In a nutshell, our methodology w.r.t. a given formal specification is applied to discover two classes of important, semantic flaws in Web sites. The first class consists of correctness errors (forbidden information that occurs in the Web site), while the second class consist of completeness errors (missing and/or incomplete Web pages). This is done by means of a novel rewriting-based technique, called *partial rewriting*, in which the traditional pattern matching mechanism is replaced by a suitable technique based on the *homeomorphic embedding* relation for recognizing patterns inside semistructured documents. The new prototype WebVerdi-M relies on a strictly more powerful Web verification engine written in Maude [12] which automatically derives the error symptoms of a given Web site. Thanks to the AC pattern matching supported by Maude and its metalevel features, we have significantly improved both the performance and the usability of the original system. By using SOAP messages and other Web-related standards, a Java Web client that interacts with a Web verification service has been made publicly available within the implementation.

Although there have been other recent efforts to apply formal techniques to Web site management [14,16,18,26], only few works addressed the semantic verification of Web sites before. The key idea behind WebVerdi-M is that rule-based techniques can support in a natural way not only intuitive, high level Web site specification, but also efficient Web site verification techniques. As far as we know, rewriting-based techniques have not been explored in the context of Web site verification to date. Previous rewriting-based approaches for Web site processing focus on transformation rather than verification issues, e.g. [23,7]. Our rule specification language does offer the expressiveness and computational power of functions and is simpler than formalizations of XML schemata based on tree automata often used in the literature (e.g. the regular expression types[21]).

VeriWeb [26] explores interactive, dynamic Web sites using a special browser that systematically explores all paths up to a specified depth. The user first specifies some properties by means of *SmartProfiles*, and then the verifier traverses the considered Web site to report the errors as sequences of Web operations that lead to a page which violates a property. Navigation errors and page errors can be signaled, but tests are performed only at the http-level. In [18], a declarative verification algorithm is developed which checks a particular class of integrity constraints concerning the Web site’s structure, but not the contents of a given instance of the site. In [14], a methodology to verify some semantic constraints concerning the Web site contents is proposed, which consists of using inference rules and axioms of natural semantics. The framework XLINKIT [16,30] allows one to check the consistency of distributed, heterogeneous documents as well as to fix the (possibly) inconsistent information. The specification language is a restricted form of first order logic combined with Xpath expressions [37] where no functions are allowed.

This work is organized as follows. Section 2 presents some preliminaries, and in Section 3 we briefly recall the rewriting-based, Web-site verification technique of [3]. In Section 4, we discuss the efficient implementation in Maude (by means of AC pattern matching) of one of the key ingredients of our verification engine: the *homeomorphic embedding* relation, which we use to recognize patterns within semi-structured documents. Section 5 briefly describes the service-oriented architecture of our verification prototype WebVerdi-M. Section 6, formalizes the API of the proposed Web Service. Section 7, we present an experimental evaluation of the system on a set of benchmarks which shows impressive performance (e.g. less than a second for evaluating a tree of some 30,000 nodes). Finally, Section 8 presented our conclusions.

2 Preliminaries

By \mathcal{V} we denote a countably infinite set of variables and Σ denotes a set of *function symbols* (also called *operators*), or *signature*. We consider varyadic signatures as in [13] (i.e., signatures in which symbols do not have a fixed arity).

$\tau(\Sigma, \mathcal{V})$ and $\tau(\Sigma)$ denote the *non-ground term algebra* and the *term algebra* built on $\Sigma \cup \mathcal{V}$ and Σ , respectively. Terms are viewed as labelled trees in the usual way. Given a term t , we say that t is *ground*, if no variable occurs in t . A *substitution* $\sigma \equiv \{X_1/t_1, X_2/t_2, \dots\}$ is a mapping from the set of variables \mathcal{V} into the set of terms $\tau(\Sigma, \mathcal{V})$ satisfying the following conditions: (i) $X_i \neq X_j$, whenever $i \neq j$, (ii) $X_i\sigma = t_i$, $i = 1, \dots, n$, and (iii) $X\sigma = X$, for all $X \in \mathcal{V} \setminus \{X_1, \dots, X_n\}$. An *instance* of a term t is defined as $t\sigma$, where σ is a substitution. By $Var(s)$ we denote the set of variables occurring in the syntactic object s . Syntactic equality between objects is represented by \equiv .

3 Rule-based Web site verification

In this section, we briefly recall the formal verification methodology proposed in [3], which allows us to detect forbidden/erroneous information as well as missing information in a Web site. This methodology is able to recognize and exactly locate the source of a possible discrepancy between the Web site and the properties required in the Web specification. An efficient and elegant implementation in Maude of such a methodology is described in Section 4.

We assume a Web page to be a well-formed *XML document* [36], since there are plenty of programs and online services that are able to validate XML syntax and perform link checking (e.g. [39],[35]). Since XML documents are provided with a tree-like structure, they can be straightforwardly encoded as ground Herbrand terms of a given term algebra.

3.1 The Web specification language

A Web specification is a triple (I_N, I_M, R) , where I_N and I_M are a finite set of correctness and completeness rules, and the set R contains the definition of some auxiliary functions.

The set I_N describes constraints for detecting erroneous Web pages (*correctness rules*). A correctness rule has the following syntax: $l \rightarrow error \mid C$ where l is a term, $error$ is a reserved constant, and C is a (possibly empty) finite sequence (which could contain membership tests of the form $X \in rexp$ w.r.t. a given regular language $rexp$;⁴ and/or equations/inequalities over terms). When C is empty, we simply write $l \rightarrow error$. Informally, the meaning of a correctness rule is the following: whenever (i) a “piece” of a given Web page can be “recognized” to be an instance $l\sigma$ of l , and (ii) the corresponding instantiated condition $C\sigma$ holds, then Web page p is marked as an incorrect page.

The third set of rules I_M specifies some properties for discovering incomplete/missing Web pages (*completeness rules*). A completeness rule is defined as $l \rightarrow r \langle \mathbf{q} \rangle$ where l and r are terms and $\mathbf{q} \in \{\mathbf{E}, \mathbf{A}\}$. Completeness rules of a Web specification formalize the requirement that some information must be included in all or some pages of the Web site. We use attributes $\langle \mathbf{A} \rangle$ and $\langle \mathbf{E} \rangle$ to distinguish “universal” from “existential” rules, as explained below. Right-hand sides r of completeness rules can contain functions, which are defined in R . In addition, some symbols in the right-hand sides of the rules may be marked by means of the symbol \sharp . Marking information of a given rule r is used to select the subset of the Web site in order to check the condition formalized by r . Intuitively, the interpretation of a universal rule (respectively, an existential rule) w.r.t. a Web site W is as follows: if (an instance of) l is “recognized” in W , (an instance of) the irreducible form of r must also be “recognized” in *all* (respectively, *some*) of the Web pages that embed (an instance of) the marked structure of r .

Diagnoses are carried out by running Web specifications on Web sites. The operational mechanism is based on a novel, flexible matching technique [3] that is able to “recognize” the partial structure of a term (Web template) within another and select it by computing *homeomorphic embeddings* (cf. [24]) of Web patterns within Web documents.

3.2 Homeomorphic embedding

Homeomorphic embedding relations allow us to verify whether a given XML document template is somehow “enclosed” within another one. We consider a simple embedding relation \trianglelefteq which closely resembles the notion of *simulation* [20], this relation has been widely used in a number of works about querying, transformation, and verification of semistructured data (cf. [10,9,1,19,8]).

Definition 1 (homeomorphic embedding). *The homeomorphic embedding relation $\trianglelefteq \subseteq \tau(\mathcal{T}ext \cup \mathcal{T}ag, \mathcal{V}) \times \tau(\mathcal{T}ext \cup \mathcal{T}ag, \mathcal{V})$ on XML documents templates is the least relation satisfying the rules:*

1. $X \trianglelefteq t$, for all $X \in \mathcal{V}$ and $t \in \tau(\mathcal{T}ext \cup \mathcal{T}ag, \mathcal{V})$.
2. $f(t_1, \dots, t_m) \trianglelefteq g(s_1, \dots, s_n)$ iff $f \equiv g$ and $t_i \trianglelefteq s_{\pi(i)}$, for $i = 1, \dots, m$, and some injective function $\pi : \{1, \dots, m\} \rightarrow \{1, \dots, n\}$.

⁴ Regular languages are represented by means of the usual Unix-like regular expression syntax.

Whenever $s \leq t$, we say that t embeds s (or s is embedded or “recognized” in t).

The intuition behind the above definition is that $s \leq t$ iff s can be obtained from t by striking out certain parts, in other words, the structure of the template s appears within t , which likely as specific Web data terms.

Let us illustrate Definition 1 by means of a rather intuitive example.

Example 1. Consider the following XML document templates (called s_1 and s_2 , respectively):

$$hpage(surname(Y), status(prof), name(X), teaching)$$

$$hpage(name(mario), surname(rossi), status(prof), \\ teaching(course(logic1), course(logic2)) \\ hobbies(hobby(reading), hobby(gardening)))$$

Note that $s_1 \leq s_2$, since the structure of s_1 can be recognized inside the structure of s_2 , while $s_2 \not\leq s_1$.

It is important to have an efficient implementation of *homeomorphic embedding* because it is used repeatedly during the verification process as described in the following.

3.3 Web verification methodology

Roughly speaking, the verification methodology works as follows. First, by using the homeomorphic embedding relation of Definition 1, we check whether the left-hand side l of some Web specification rule is embedded into a given page p of the considered Web site. When the embedding test $l \leq p$ succeeds, by extending the proof, we construct the biggest substitution⁵ σ for the variables in $Var(l)$, such that $l\sigma \leq p$. Then, depending on the nature of the Web specification rule (correction or completeness rule), it is as follows:

- (**Correction rule**) evaluating the condition of the rule (instantiated by σ); a correctness error is signalled in the case when the error condition is fulfilled.
- (**Completeness rule**) by a new homeomorphic embedding test, checking whether the right-hand side of the rule (instantiated by σ) is recognized in some page of the considered Web site. Otherwise, a completeness error is signalled. Moreover, from the incompleteness symptom computed so far, a fixpoint computation is started in order to discover further missing information, which may involve the execution of other completeness rules.

⁵ The substitution σ is easily obtained by composing the bindings X/t , which can be recursively gathered during the *homeomorphic embedding* test $X \leq t$, for $X \in l$ and $t \in p$.

4 Verifying Web sites using Maude

Maude is a high-performance reflective language supporting both equational and rewriting logic programming, which is particularly suitable for developing domain-specific applications [33,15]. In addition, the Maude language is not only intended for system prototyping, but it has to be considered as a real programming language with competitive performance. In the rest of the section, we recall some of the most important features of the Maude language which we have conveniently exploited for the optimized implementation of our Web site verification engine.

Equational attributes. Let us describe how we model (part of) the internal representation of XML documents in our system. The chosen representation slightly modifies the data structure provided by the Haskell HXML Library [17] by adding commutativity to the standard XML tree-like data representation. In other words, in our setting, the order of the children of a tree node is not relevant: e.g., $f(a, b)$ is “equivalent” to $f(b, a)$.

```
fmod TREE-XML is
sort XMLNode .
op RTNode : -> XMLNode .           -- Root (doc) information item
op ELNode _ _ : String AttList -> XMLNode . -- Element information item
op TXNode _ : String -> XMLNode .     -- Text information items
--- ... definitions of the other XMLNode types omitted ...
sorts XMLTreeList XMLTreeSeq XMLTree .
op Tree ( _ ) _ : XMLNode XMLTreeList - > XMLTree .
subsort XMLTree < XMLTreeSeq .
op _,_ : XMLTreeSeq XMLTreeSeq -> XMLTreeSeq [comm assoc id:null] .
op null : -> XMLTreeSeq .
op [_] : XMLTreeSeq -> XMLTreeList .
op [] : -> XMLTreeList .
endfm
```

In the previous module, the XMLTreeSeq constructor `_,_` is given the equational attributes `comm assoc id:null`, which allow us to get rid of parentheses and disregard the ordering among XML nodes within the list. The significance of this optimization will be clear when we consider rewriting XML trees with AC pattern matching.

AC pattern matching. The evaluation mechanism of Maude is based on rewriting modulo an equational theory E (i.e. a set of equational axioms), which is accomplished by performing *pattern matching modulo* the equational theory E . More precisely, given an equational theory E , a term t and a term u , we say that t *matches* u modulo E (or that t *E -matches* u) if there is a substitution σ such that $\sigma(t) =_E u$, that is, $\sigma(t)$ and u are equal modulo the equational theory E . When E contains axioms for associativity and commutativity of operators, we talk about *AC pattern matching*. AC pattern matching is a powerful matching mechanism, which we employ to inspect and extract the partial structure of

a term. That is, we use it directly to implement the notion of homeomorphic embedding of Definition 1.

Metaprogramming. Maude is based on rewriting logic [25], which is reflective in a precise mathematical way. In other words, there is a finitely presented rewrite theory \mathcal{U} that is universal in the sense that we can represent in \mathcal{U} (as a data) any finitely presented rewrite theory \mathcal{R} (including \mathcal{U} itself), and then mimic in \mathcal{U} the behavior of \mathcal{R} . We have used the metaprogramming capabilities of Maude to implement the semantics of correctness as well as completeness rules (e.g. implementing the homeomorphic embedding algorithm, evaluating conditions of conditional rules, etc.). Namely, during the partial rewriting process, functional modules are dynamically created and run by using the meta-reduction facilities of the language.

Now we are ready to explain how we implemented the homeomorphic embedding relation of Section 3.2, by exploiting the aforementioned Maude high-level features.

4.1 Homeomorphic embedding implementation.

Let us consider two XML document templates l and p . The critical point of our methodology is to (i) discover whether $l \trianglelefteq p$ (i.e. l is embedded into p); (ii) find the substitution σ such that $l\sigma$ is the instance of l recognized inside p , whenever $l \trianglelefteq p$.

Given l and p , our proposed solution can be summarized as follows. By using Maude metalevel features, we first dynamically build a module M that contains a single rule of the form

$$\text{eq } l = \text{sub}(\text{"X}_1\text{"/X}_1), \dots, \text{sub}(\text{"X}_n\text{"/X}_n), \quad \text{X}_i \in \text{Var}(l), i = 1, \dots, n,$$

where **sub** is an associative operator used to record the substitution σ that we want to compute. Next, we try to reduce the XML template p by using such a rule. Since l and p are internally represented by means of the binary constructor `_,_` that is given the equational attributes `comm assoc id:null` (see Section 4), the execution of module M on p essentially boils down to computing an AC-matcher between l and p . Moreover, since AC pattern matching directly implements the homeomorphic embedding relation. The execution of M corresponds to finding all the homeomorphic embeddings of l into p (recall that the set of AC matchers of two compatible terms is not generally a singleton). Additionally, as a side effect of the execution of M , we obtain the computed substitution σ for free as the sequence of bindings for the variables X_i , $i = 1, \dots, n$ which occur in the instantiated rhs

$$\text{sub}(\text{"X}_1\text{"/X}_1)\sigma, \dots, \text{sub}(\text{"X}_n\text{"/X}_n)\sigma, \quad \text{X}_i \in \text{Var}(l), i = 1, \dots, n,$$

of the dynamic rule after the partial rewriting step.

Example 2. Consider again the XML document templates s_1 and s_2 of Example 1. We build the dynamic module M containing the rule

$$\text{op } \text{hpage}(\text{surname}(Y), \text{status}(\text{prof}), \text{name}(X), \text{teaching}) = \text{sub}(\text{"Y"/Y}, \text{sub}(\text{"X"/X}) .$$

Since $s_1 \leq s_2$, there exists an AC-match between s_1 and s_2 and, hence, the result of executing M against the (ground) XML document template s_2 is the computed substitution: $\text{sub}("Y"/\text{rossi}), \text{sub}("X"/\text{mario})$.

5 Prototype implementation

The application has been structured as a SOAP Web Service [38]. The main reason behind this choice is the advantage of using the Service Oriented Architecture paradigm [40,22]. In this paradigm services are distributed, autonomous, and independent. They are realized using standard protocols, in order to build networks of collaborating applications. This style of architecture allows one to reuse at the macro level (service), rather than micro level (object).

WebVerdi-M as a service-oriented architecture allows to access the core verification engine Verdi-M as a reusable entity. This implementation is public available at <http://www.dsic.upv.es/users/elp/webverdi-m>.

WebVerdi-M can be divided into two layers: front-end and back-end. The back-end layer provides web services to support the front-end layer. This architecture allow clients on the network to invoke the Web service functionality through the available interfaces.

The tool consists of the following components: Web service WebVerdiService, Web client WebVerdiClient, core engine Verdi-M, XML API, and database DB. Figure 1 illustrates the overall architecture of the system.

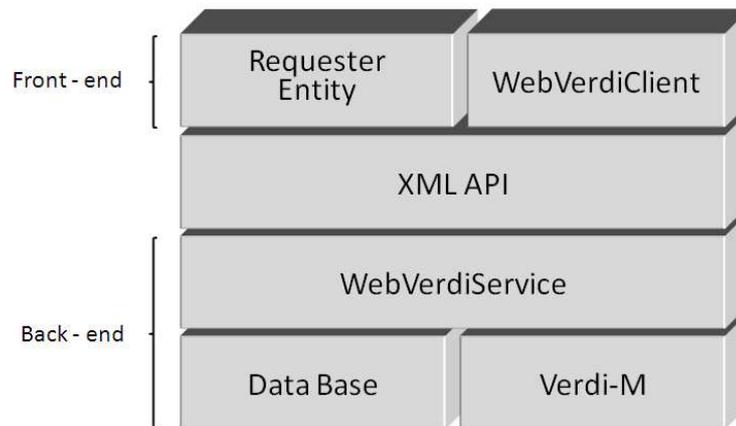


Fig. 1. Components of WebVerdi-M

5.1 WebVerdiService

The sever has been structured as a Web Service and its control parts is implemented in Java 1.4 [27], in order to use the large number of implementations of

the Web Service standards available as the TriActive JDO persistence package [34]. Persistence is used to store both web pages and specification locally to the sever in a MySQL database [29].

TriActive JDO is an open source implementation of Sun's JDO specification [28], designed to support transparent persistence using any JDBC-compliant database. TriActive JDO allows to generate schema, meaning it takes user-written Java classes and automates the tasks required to transparently persist objects to a database.

All the elaboration is performed on the server side by modules written in Maude [32,31]. This flavor of functional programming has been used because of its meta-level functionality, which together with the properties of associativity and commutativity, allows a simple and elegant implementation of the simulation algorithm.

The web service exports six operations that are network-accessible through standardized XML messaging. These operations are: store a Web site, remove a Web site, retrieve a Web site, add Web page to a Web site, check correctness, and check completeness. The Web service acts as a single access point to the core engine Verdi-M which implements the Web verification methodology in Maude. Following the standards, the architecture is also platform and language independent so as to be accessible via scripting environment as well as via client applications across multiple platforms.

5.2 XML API

In order for successful communications to occur, both the WebVerdiService and WebVerdiClient (or any user) must agree to a common format for the messages being delivered so that they can be properly interpreted at each end. The WebVerdiService Web service is developed by defining an API (see Section 6) that encompasses the executable library of the core engine. This is achieved by making use of Oracle JDeveloper, including the generation of WSDL for making the API available. The OC4J Server (the web server integrated in Oracle JDeveloper) handles all procedures common to Web service development. Synthesized error symptoms are also encoded as XML documents in order to be transferred from the WebVerdiService Web service to client applications as an XML response by means of the SOAP protocol.

5.3 Verdi-M

Verdi-M is the most important part of the tool. Here is where the verification methodology is implemented. This component is implemented in Maude language and is independent of the other system components. The module of engine are described in [2]. These modules are invoked as separate processes when needed by the Java main thread. These modules, using the functionalities of the core. Overall the flow execution of the elaboration is Java controlled.

5.4 WebVerdiClient

The client consists of a Java graphical interface which allows to use the functionalities offered by the Web Server. The client is using the API specified in Section 6 to interact with it through a network connection and uses the Java Web Start functionalities to execute. The main goal was to provide an *intuitive* and *friendly* interface for the user.

WebVerdiClient is provided with a versatile, new graphical interface that offers three complementary views for both the specification rules and the pages of the considered Web site: the first one is based on the typical idea of accessing contents by using folders trees and is particularly useful for beginners; the second one is based on XML, and the third one is based on term algebra syntax. The tool provides all translations among the three views.

A snapshot of WebVerdiClient is shown in Figure 2.

Any other client using the API of the Web Server could be used.

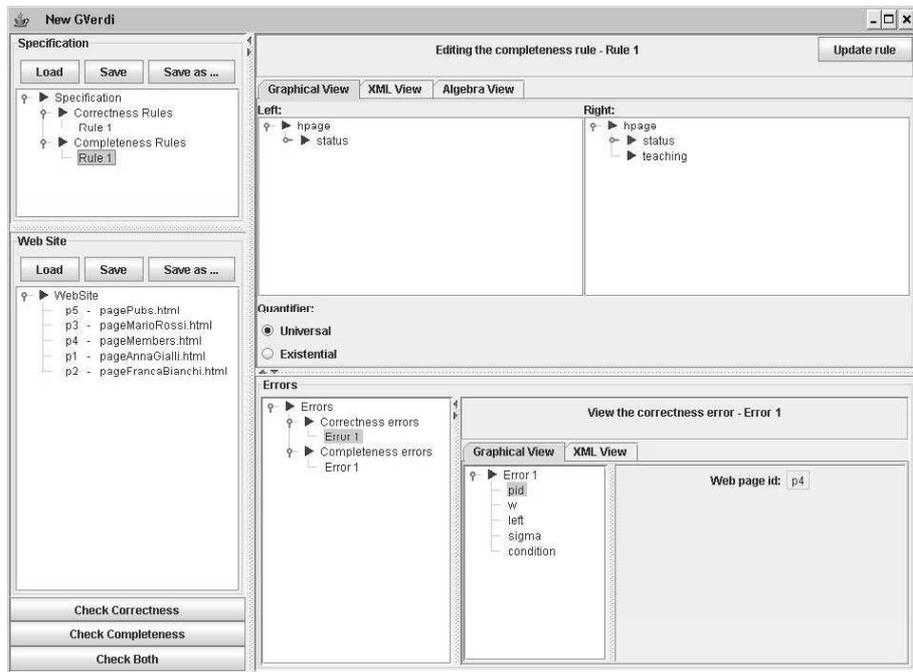


Fig. 2. WebVerdiClient Snapshot

5.5 DB

The WebVerdiService Web service needs to transmit abundant XML data over the Web to and from client applications. The common behavior of a user when using

the tool is to modify the default rules provided for the Web specification and then verify a particular Web site. After modifying the Web specification, it would be necessary to send back to the service the considered specification as well as the whole Web site to verify. After the application invokes the WebVerdiService Web service with these two elements, synthesized errors are progressively generated and transferred to the client application. The standard Web service architecture requires client applications to wait until all data are received and then errors are sent, which could cause significant time lags in the application. In order to avoid this overhead and to provide better performance to the user, we use a local *MySQL* data base where the Web site and Web errors are temporarily stored at the server side.

6 API

This section summarizes the types of methods and the specific message exchange patterns that are considered for interacting with WebVerdi-M.

We considerer the data representation and methods separately.

6.1 Data Representation

Web Site. A website consist of a set of pages, which are represented as a triple (id,name,data) where:

id: is the page identifier;

name: is the name of the page;

data: is the XML content of the page.

An example of page is given in Example 3.

Example 3. Example of Web page:

```
<page>
  <id>p1</id>
  <name>biblio.htm</name>
  <data>
    <biblioteca>
      <libro ndoc="99231">
        <isbn>8437607000</isbn>
        <autor>Rojas, Fernando de</autor>
        <titulo>La Celestina</titulo>
      </libro>
      <libro ndoc="158290">
        <isbn>8403870485</isbn>
        <autor>Homero</autor>
        <titulo>Iliada</titulo>
      </libro>
      <libro ndoc="181227">
        <isbn>8466401040</isbn>
```

```

        <autor>Kafka, Franz</autor>
        <titulo>La metamorfosi</titulo>
    </libro>
</biblioteca>
</data>
</page>

```

A Web Site is represented by a collection of web pages. The XML encoding of a Web Site is shown in Example 4.

Example 4. Example of web site:

```

<webSite>
  <page> ... </page>
  <page> ... </page>
  <page> ... </page>
  ...
  <page> ... </page>
</webSite >

```

Rules. In our methodologies, there are two different kinds of rules, which are handled differently. Different XML representations are therefore needed.

Correctness rules A correctness rule is defined by the triple (l,r,C), where:

l: is the left-hand side of the rule;
r: is the right-hand side of the rule;
C: is the condition (if any).

The XML representation of a correctness rule is as follow.

```

<ruleCorrectness Name="...">
  <left> ... </left>
  <right> ... </right>
  <condition> ... </condition>
</ruleCorrectness>

```

Completeness rules A completeness rule is encoded by the triple (l,r,q), where:

l: is the left-hand side of the rule;
r: is the right-hand side of the rule;
q is the logical quantifier, which can be either E (Existential) o A (Universal).

An attribute is added to identify a particular rule inside the Web specification. The XML representation of a completeness rule is as follow

```

<ruleCompleteness Name="Rule 1">
  <left>
    <atrib>
      <f>
        <atrib>
          <g>

```

```

        <atrib>X</atrib>
      </g>
    </atrib>
  </f>
</atrib>
</left>
<right>
  <atrib Mark=true>
    <h>
      <atrib Mark=true>
        <g>
          <atrib>X</atrib>
        </g>
      </atrib>
    </h>
  </atrib>
</righth>
<quantifier>E</quantifier>
</ruleCompletness>

```

A Web specification is a collection of completeness and/or correctness rules. Its XML representation is drafted in Example 5.

Example 5. Specification example:

```

<specification>
  <ruleCorrectness Name=xxx> ... </ruleCorrectness>
  <ruleCompletness Name=xxx> ... </ruleCompletness>
  <ruleCompletness Name=xxx> ... </ruleCompletness>
  ...
  <ruleCorrectness Name=xxx> ... </ruleCorrectness>
</specification>

```

Errors. Correction errors are given by a 4-pla (pid,w,l,sigma,C), where:

- pid:** is the identifier of the page which contains errors;
- w:** is the position within the page where the error is located; this position is defined by a vector of integers like [1, 4, 5]; in XML this is written "1.4.5";
- l:** the left-hand side of the rule which produces the error;
- sigma:** the substitution(s) in *l* which produce(s) the error;
- C:** the condition of the rule which produces the error.

Example 6. Correctness error example:

```

<errorCorrectness>
  <pid>p1</pid>
  <w>1.2</w>
  <l>
    <autor>X</autor>
  </l>

```

```

<sigma>
  <sust>
    <var>X</var>
    <value>Rojas, Fernando de</value>
  </sust>
</sigma>
<condition>X=Rojas, Fernando de</condition>
</errorCorrectness>

```

The completeness errors representation depends on the type of error: missing page error or universal/existential error.

- Missing Page Errors are defined by the triple (r,W,sigma):
 - r**: Rule which generates the error;
 - W**: Web Site;
 - sigma**: the substitution(s) which produce(s) the error.
- Universal/Existential Errors are defined by the triple (r,P,sigma):
 - r**: Rule which generates the error;
 - P**: The set of identifiers of pages which do not comply with the rule;
 - sigma**: the substitution(s) which produce(s) the error.

An attribute is needed to distinguish among different types of errors. The values of this attribute can be either M (Missing Page), A (Universal), or E (Existential). Its XML representation is given in Example 7.

Example 7. Completeness error example:

```

<errorCompletness>
  <r> ... </r>
  <pages>
    <pid>p1</pid>
    <pid> ... </pid>
    <pid>pn</pid>
  </pages>
  <sigma> ... </sigma>
  <type> ... </type>
</errorCompletness>

```

A collection of errors is a set of completeness or correctness errors which is represented in XML as follows.

Example 8. Error collection example:

```

<collectionErrors>
  <errorCorrectness> ... </errorCorrectness>
  ...
  <errorCorrectness> ... </errorCorrectness>
</collectionErrors>

...

```

```

<collectionErrors>
  <errorCompleteness> ... </errorCompleteness>
  ...
  <errorCompleteness> ... </errorCompleteness>
</collectionErrors>

```

Actions. An action is the primitive to repair the website. There are four different primitives:

change (pid,w,t): Changes the subterm in the position w of the page pid with the term t;

insert(pid,w,t): Adds term t in the position w of the page pid;

delete(pid,t): Deletes the term t from the page pid;

add(p, idWS): Adds the page p to the website idWS.

Example 9. Action example:

```

<action>
  <type>change</type>
  <pid>p1</pid>
  <w>1.2</w>
  <t>
    <autor>Perez, Pepito</autor>
  </t>
</action>

```

PairErrorAction. Each kind of error is handled by an specific action. For instance, in the case of an error generated by a Universal rule, the repairing action applies to all the pages that do not fulfill the rule.

Its XML representation is as in Example 10.

Example 10. PairErrorAction example:

```

<pairErrorAction>
  <errorCorrectness> ... </errorCorrectness>
  <action> ... </action>
</pairErrorAction>

...

<pairErrorAction>
  <errorCompleteness> ... </errorCompleteness>
  <action> ... </action>
</pairErrorAction>

```

6.2 Methods

To optimize the data transfer, some methods are provided as defined below (storeWebSite, recoveryPage, recoveryWebSite and removeWebSite).

Methods Descriptions.

<p>storeWebSite (WebSite)</p> <p>Stores a Web Site in the local client server.</p> <p>Input: WebSite: The Web Site to be stored.</p> <p>Output: The identifier of the Web Site stored in the local server.</p>
<p>retrievePage(idP, idWS)</p> <p>Loads a page of the specified Web Site.</p> <p>Input: idP: The identifier of the page. idWS: The identifier of the Web Site stored in the local server.</p> <p>Output: The retrieved page of the website, if the identifiers exists.</p>
<p>retrieveWebSite(idWS)</p> <p>Loads the Web Site stored in the local server.</p> <p>Input: idWS: the identifiers of the website as stored in the local server.</p> <p>Output: The Web Site, if the identifiers exist.</p>
<p>removeWebSite(idWS)</p> <p>Deletes the Web Site stored in the local server.</p> <p>Input: idWS: the Web Site Identifier stored in the local server.</p> <p>Output: True if the website has been successfully deleted; False otherwise.</p>
<p>checkCorrectness(idWS SPEC)</p> <p>Returns a collection of correctness errors of a Web Site, w.r.t. the given specification.</p> <p>Input: idWS: The identifier of the Web Site as stored in the local server. SPEC: The XML representation of the Web specification (see 6.1).</p> <p>Output:</p>

An XML encoding collection of a correctness errors (see 6.1).
checkCompleteness(idWS SPEC)
<p>Returns a collection of completeness errors of a Web Site, w.r.t. the given specification.</p> <p>Input: idWS: The identifier of the Web Site as stored in the local server. SPEC: The XML representation of the Web specification (see 6.1).</p> <p>Output: An XML encoding collection of a completeness errors (see 6.1).</p>
fixErrorCorrectnessByDelete (errorCorrecness idWS)
<p>Repairs a completeness error by a deletion action.</p> <p>Input: errorCorrecness: Error to be fixed. idWS: The identifier of the Web Site as stored in the local server.</p> <p>Output: True in case the error was successfully deleted; False otherwise.</p>
changeCS (errorCorrecness SPEC idWS)
<p>Repairs automatically the correctness error by using of a constraint solver.</p> <p>Input: errorCorrecness: Error to be fixed. SPEC: The XML representation of the web specification (see 6.1). idWS: The identifier of the Web Site as stored in the local server.</p> <p>Output: True in case the error was successfully solved; False otherwise.</p>
fixErrorCompleteness (pairErrorAction SPEC idWS)
<p>Repairs a completeness error.</p> <p>Input: pairErrorAction: Pair error-action (see 6.1). SPEC: The XML representation of the Web specification (see 6.1). idWS: The identifier of the Web Site as stored in the local server.</p> <p>Output: True in case the error can be fixed; False otherwise.</p>

fixErrorCompletenessByStrategyM (idWS pairEA_Set)
Repairs a completeness error using a Minimal strategy.
Input: idWS: The identifier of the Web Site as stored in the local server. pairEA_Set: The set of couples <error,action> (see 6.1).
Output: True in case the error can be fixed; False otherwise.
fixErrorCompletenessByStrategyMNO (idWS pairEA_Set)
Repairs a completeness error using a Minimal non-overlapping strategy.
Input: idWS: The identifier of the Web Site as stored in the local server. pairEA_Set: The set of couples <error,action> (see 6.1).
Output: True in case the error can be fixed; False otherwise.

7 Experimental evaluation

In order to evaluate the usefulness of our approach in a realistic scenario (that is, for sites whose data volume exceeds toy sizes), we have benchmarked our system by using several correctness as well as completeness rules of different complexity for a number of XML documents randomly generated by using the XML documents generator `xmlgen` available within the XMark project [11]. The tool `xmlgen` is able to produce a set of XML data, each of which is intended to challenge a particular primitive of XML processors or storage engines by using different scale factors.

Table 2 shows some of the results we obtained for the simulation of three different Web specifications *WS1*, *WS2* and *WS3* in five different, randomly generated XML documents. Specifically, we tuned the generator for scaling factors from 0.01 to 0.1 to match an XML document whose size ranges from 1Mb –corresponding to an XML tree of about 31000 nodes– to 10Mb –corresponding to an XML tree of about 302000 nodes– (for an exhaustive evaluation, including comparison with related systems, please refer to <http://www.dsic.upv.es/users/elp/webverdi-m/>).

Both Web specifications *WS1* and *WS2* aim at checking the verification power of our tool regarding data correctness, and thus include only correctness

rules. The specification rules of *WS2* contain more complex and more demanding constraints than the ones formalized in *WS1*, with involved error patterns to match, and conditional rules with a number of membership tests and functions evaluation. The Web specification *WS3* aims at checking the completeness of the randomly generated XML documents. In this case, some critical completeness rules have been formalized which recognize a significant amount of missing information.

Size	Nodes	Scale factor	Time		
			<i>WS1</i>	<i>WS2</i>	<i>WS3</i>
1 Mb	30,985	0.01	0.930 s	0.969 s	165.578 s
3 Mb	90,528	0.03	2.604 s	2.842 s	1768.747 s
5 Mb	150,528	0.05	5.975 s	5.949 s	4712.157 s
8 Mb	241,824	0.08	8.608 s	9.422 s	12503.454 s
10 Mb	301,656	0.10	12.458 s	12.642 s	21208.494 s

Table 2. Verdi-M Benchmarks

The results shown in Table 2 were obtained on a personal computer equipped with 1Gb of RAM memory, 40Gb hard disk and a Pentium Centrino CPU clocked at 1.75 GHz running Ubuntu Linux 5.10.

Let us briefly comment our results. Regarding the verification of correctness, the implementation is extremely time efficient, with elapsed times scaling lineary. Table 2 shows that the execution times are small even for very large documents (e.g. running the correctness rules of Web specification *WS1* over a 10Mb XML document with 302000 nodes takes less than 13 seconds). Concerning the completeness verification, the fixpoint computation which is involved in the evaluation of the completeness rules typically burdens the expected performance (see [3]), and we are currently able to process efficiently XML documents whose size is not bigger than 1Mb (running the completeness rules of Web specification *WS3* over a 1Mb XML document with 31000 nodes takes less than 3 minutes).

Finally, we want to point out that the current Maude implementation of the verification system supersedes and greatly improves our preliminary system, called *GVerdi*[3,6], that was only able to manage correctness for small XML repositories (of about 1Mb) within a reasonable time. We are currently working on further improving the performance of our system.

8 Conclusion

In the literature on Web management, Web sites verification has mainly a syntactic focus with a particular concern for the accessibility and usability perspective [4,5]. This paper can be seen as a step forward towards the formal, semantic verification of Web sites using rule-based technology. First we present an efficient and innovative implementation in Maude –a high-performance reflective

functional language— of the rewriting-based, Web verification methodology of [3]. This methodology deals with semantic flaws that are not addressed by classical tools. The framework comes with a language for defining correctness and completeness conditions on Web sites. Then, our rewriting-based verification technique is able to recognize forbidden/incorrect patterns and incomplete/missing Web pages by means of a novel rewriting-based technique, called *partial rewriting*.

In this work, first we exploit Maude’s capabilities which are particularly suitable for our implementation, such as associative commutative pattern matching and metaprogramming. We can thus provide WebVerdi-M with a powerful Web verification engine. We have done a comparison of run times of Verdi-M core engine and shown the resulting impressive performance (e.g. less than a second for evaluating a tree of some 30,000 nodes).

Then, we have proposed a service-oriented architecture which makes the Web verification capabilities of the system easily accessible to internet requestors. The resulting prototype WebVerdi-M is publicly available together with a set of examples and its XML API.

In order to make possible technological transfer to industry it is necessary to have tools that are able to give prompt answers on real size examples, as we have shown by our scalable benchmarks. Another important factor, is to reduce the cost of learning to the user. For this reason we have developed a friendly innovative interface for our system.

References

1. S. Abiteboul, P. Buneman, and D. Suciu. *Data on the Web. From Relations to Semistructured Data and XML*. Morgan Kaufmann, 2000.
2. M. Alpuente, D. Ballis, S. Escobar, M. Falaschi, P. Ojeda, and D. Romero. Un motor algebraico para la verificacin de sistemas Web en Gverdi. Technical Report DSIC-II/02/07, DSIC, UPV, 2007.
3. M. Alpuente, D. Ballis, and M. Falaschi. Automated Verification of Web Sites Using Partial Rewriting. *Software Tools for Technology Transfer*, 8:565–585, 2006.
4. M. Alpuente, S. Escobar, and M. Falaschi (Eds.). *Automated Specification and Verification of Web Sites, 1st Int’l Workshop WWV’05*, volume 157(2). Elsevier, 2006.
5. M. Alpuente, S. Escobar, and M. Falaschi (Eds.). *Automated Specification and Verification of Web Systems, 2nd Int’l Workshop WWV’06*. IEEE Computer Society Press, 2007.
6. D. Ballis and J. García Vivó. A Rule-based System for Web Site Verification. In *Proc. of 1st Int’l Workshop on Automated Specification and Verification of Web Sites (WWV’05)*. ENTCS, Elsevier, 2005. To appear.
7. I. D. Baxter, F. Ricca, and P. Tonella. Web Application Transformations based on Rewrite Rules. *Information and Software Technology*, 44(13), 2002.
8. E. Bertino, M. Mesiti, and G. Guerrin. A Matching Algorithm for Measuring the Structural Similarity between an XML Document and a DTD and its Applications. *Information Systems*, 29(1):23–46, 2004.

9. F. Bry and S. Schaffert. Towards a Declarative Query and Transformation Language for XML and Semistructured Data: Simulation Unification. In *Proc. of the Int'l Conference on Logic Programming (ICLP'02)*, volume 2401 of *LNCS*. Springer-Verlag, 2002.
10. F. Bry and S. Schaffert. The XML Query Language Xcerpt: Design Principles, Examples, and Semantics. Technical report, 2002. Available at: <http://www.xcerpt.org>.
11. Centrum voor Wiskunde en Informatica. XMark – an XML Benchmark Project, 2001. Available at: <http://monetdb.cwi.nl/xml/>.
12. Manuel Clavel, Francisco Durán, Steven Eker, Patrick Lincoln, Narciso Martí-Oliet, José Meseguer, and Carolyn Talcott. The maude 2.0 system. In Robert Nieuwenhuis, editor, *Rewriting Techniques and Applications (RTA 2003)*, number 2706 in *Lecture Notes in Computer Science*, pages 76–87. Springer-Verlag, 2003.
13. N. Dershowitz and D. Plaisted. Rewriting. *Handbook of Automated Reasoning*, 1:535–610, 2001.
14. T. Despeyroux and B. Trousse. Semantic Verification of Web Sites Using Natural Semantics. In *Proc. of 6th Conference on Content-Based Multimedia Information Access (RIAO'00)*, 2000.
15. S. Eker, J. Meseguer, and A. Sridharanarayanan. The Maude LTL model checker and its implementation. In *Model Checking Software: Proc. 10th Intl. SPIN Workshop*, volume 2648 of *LNCS*, pages 230–234. Springer, 2003.
16. E. Ellmer, W. Emmerich, A. Finkelstein, and C. Nentwich. Flexible Consistency Checking. *ACM Transaction on Software Engineering*, 12(1):28–63, 2003.
17. Joe English. The HXML Haskell Library, 2002. Available at: <http://www.flightlab.com/joe/hxml/>.
18. M. Fernandez, D. Florescu, A. Levy, and D. Suciu. Verifying Integrity Constraints on Web Sites. In *Proc. of Sixteenth International Joint Conference on Artificial Intelligence (IJCAI'99)*, volume 2, pages 614–619. Morgan Kaufmann, 1999.
19. M. F. Fernandez and D. Suciu. Optimizing Regular Path Expressions Using Graph Schemas. In *Proc. of Int'l Conference on Data Engineering (ICDE'98)*, pages 14–23, 1998.
20. M. R. Henzinger, T. A. Henzinger, and P. W. Kopke. Computing Simulations on Finite and Infinite Graphs. In *IEEE Symposium on Foundations of Computer Science*, pages 453–462, 1995.
21. H. Hosoya and B. Pierce. Regular Expressions Pattern Matching for XML. In *Proc. of 25th ACM SIGPLAN-SIGACT Int'l Symp. POPL*, pages 67–80. ACM, 2001.
22. IBM. Service Oriented Architecture, 2007. Available at: <http://www-306.ibm.com/software/solutions/soa>.
23. C. Kirchner, Z. Qian, P. K. Singh, and J. Stuber. Xemantics: a Rewriting Calculus-Based Semantics of XSLT. Rapport de recherche A01-R-386, LORIA, 2001.
24. M. Leuschel. Homeomorphic Embedding for Online Termination of Symbolic Methods. In T. Æ. Mogensen, D. A. Schmidt, and I. H. Sudborough, editors, *The Essence of Computation*, volume 2566 of *Lecture Notes in Computer Science*, pages 379–403. Springer, 2002.
25. N. Martí-Oliet and J. Meseguer. Rewriting Logic: Roadmap and Bibliography. *Theoretical Computer Science*, 285(2):121–154, 2002.
26. B. Michael, F. Juliana, and G. Patrice. Veriweb: automatically testing dynamic web sites. In *Proc. of 11th Int'l WWW Conference*. ENTCS, Elsevier, 2002.
27. Sun Microsystems. Java™ 2 SDK, Standard Edition Documentation, Version 1.4.2, 2003. Available at: <http://java.sun.com>.

28. Sun Microsystems. JSR 12: Java™ Data Objects (JDO) Specification, 2006. Available at: <http://www.jcp.org/en/jsr/detail?id=12>.
29. MySQLAB. MySQL, 2007. Available at: <http://www.mysql.com>.
30. C. Nentwich, W. Emmerich, and A. Finkelstein. Consistency Management with Repair Actions. In *Proc. of the 25th International Conference on Software Engineering (ICSE'03)*. IEEE Computer Society, 2003.
31. Department of Computer Science Univeristy of Illinois at Urbana-Champaign. Maude manual and examples. Available at: <http://maude.cs.uiuc.edu>.
32. Department of Computer Science Univeristy of Illinois at Urbana-Champaign. Maude primer and examples. Available at: <http://maude.cs.uiuc.edu>.
33. J. Meseguer S. Escobar, C. Meadows. A Rewriting-Based Inference System for the NRL Protocol Analyzer and its Meta-Logical Properties. *Theoretical Computer Science*, 367(1-2):162–202, 2006.
34. SourceForge. TriActive JDO, 2005. Available at: <http://tjdo.sourceforge.net>.
35. Typke und Wicke GbR. Validate/Check XML. Available at: <http://www.xmlvalidation.com/>.
36. World Wide Web Consortium (W3C). Extensible Markup Language (XML) 1.0, second edition, 1999. Available at: <http://www.w3.org>.
37. World Wide Web Consortium (W3C). XML Path Language (XPath), 1999. Available at: <http://www.w3.org>.
38. World Wide Web Consortium (W3C). SOAP Version 1.2, 2003. Available at: <http://www.w3.org>.
39. World Wide Web Consortium (W3C). Markup Validation Service, 2005. Available at: <http://validator.w3.org/>.
40. Wikipedia. Service Oriented Architecture, 2007. Available at: <http://en.wikipedia.org>.