

## Rule-based verification of Web sites

M. Alpuente · D. Ballis · M. Falaschi

Published online: 26 July 2006  
© Springer-Verlag 2006

**Abstract** In this paper, we develop a framework for the automated verification of Web sites, which can be used to specify integrity conditions for a given Web site, and then automatically check whether these conditions are fulfilled. First, we provide a rewriting-based, formal specification language which allows us to define syntactic as well as semantic properties of the Web site. Then, we formalize a verification technique which detects both incorrect/forbidden patterns as well as lack of information, that is, incomplete/missing Web pages inside the Web site. Useful information is gathered during the verification process which can be used to repair the Web site. Our methodology is based on a novel rewriting-based technique, called *partial rewriting*, in which the traditional pattern matching mechanism is replaced by tree *simulation*, a suitable technique for recognizing patterns inside semistructured documents. The framework has been implemented in the prototype GVERDI, which is publicly available.

### 1 Introduction

The increasing complexity of Web sites has turned their design and construction into a challenging problem. Systematic, formal approaches can bring many benefits to Web site construction, giving support for automated Web site verification. This paper presents an approach to Web site specification and verification based on rewriting-like machinery. We use rewriting-based technology both to specify the integrity conditions and to formalize a verification technique which obtains the requirements not fulfilled by the Web site, and then is able to diagnose errors by finding out incorrect/forbidden patterns and missing/incomplete pages.

Although the management of Web sites has received significant attention in recent years [10,16,21], few works address the semantic verification of Web sites. In [21], 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 [16], 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 [19,32] 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. With respect to the correctness of Web applications, a symbolic model-checking approach is formalized in [17], which constructs a finite states model of the system in the model checker input language, and then checks the considered properties which are expressed in CTL logic. For a comprehensive

---

M. Alpuente (✉)  
DSIC, Universidad Politécnica de Valencia,  
Camino de Vera s/n, Apdo. 22012, 46071 Valencia, Spain  
e-mail: alpuente@dsic.upv.es

D. Ballis  
Dip. Matematica e Informatica,  
Via delle Scienze 206, 33100 Udine, Italy  
e-mail: demis@dimi.uniud.it

M. Falaschi  
Dip. di Scienze Matematiche e Informatiche,  
Pian dei Mantellini 44, 53100 Siena, Italy  
e-mail: moreno.falaschi@unisi.it

survey about the general problem of checking constraints between multiple documents, we refer to [20, 18, 13].

Our idea in this paper is that term rewriting 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. We only know of two rewriting-based approaches for Web site processing, but they focus on transformation rather than verification issues: a term rewriting implementation is provided in [29] for (a fragment of) XSLT, the rule-based language designed by W3C for the transformation of XML documents, whereas rewrite rules are used in [7] to perform HTML transformations with the aim of improving Web applications by cleaning up syntax, reorganizing frames, or updating to new standards. 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 such as, e.g., the regular expression types [26].

#### Our contribution

We first provide a rewriting-based, formal specification language which allows us to define conditions on both the structure and the contents of Web sites in a simple and concise way. For instance, it allows us to recognize erroneous information inside the Web site and, additionally, to enforce that some information is available at a given Web page, including the adequate links between pages or even the existence of the Web pages themselves. In our formalism, web pages (HTML/XML documents) are modeled as Herbrand terms, and, consequently, Web sites are finite sets of terms. Then, we formalize a verification technique in which a Web site is checked w.r.t. a given Web specification in order to detect incorrect data and incomplete and/or missing Web pages. Moreover, by analyzing the error symptoms gathered during the verification process, we are also able to (1) exactly locate the incorrect/forbidden information and (2) to find out the missing information which the user should provide to repair the Web site. Since reasoning on the Web calls for formal methods specifically fitting the Web context, we combine a standard regular expressions methodology with a novel, rewriting-based technique called *partial rewriting*, in which the traditional pattern matching mechanism is replaced by tree *simulation* [24] in order to provide a suitable mechanism for recognizing patterns inside semistructured documents. The notion of simulation has been already used for dealing with semistructured data in a number of query and transfor-

mation languages [10, 12, 22, 14]. The reason is twofold: on the one hand, it provides a powerful method to extract information from semistructured data; on the other hand, there exist efficient algorithms for computing simulations [24]. To assess the feasibility and efficiency of our approach, we have implemented the prototype system GVERDI (Graphical VERification and Rewriting for Debugging Internet sites), which is based on the verification methodology that we propose and is publicly available online. Following the “tolerant” approach of XLINKIT [19, 32], we do not force the immediate repairing of the Web site, but simply provide the diagnosis information that enables document owners to decide on further actions. Some of the results in this work have been included in Ballis’ doctoral thesis [5].

#### Plan of the paper

The rest of the paper is organized as follows. Section 2 summarizes some preliminary definitions and notations. In Sect. 3, we formulate a simple method for translating XHTML/XML documents into Herbrand terms. Section 4 is devoted to formalize the specification language, whereas Sect. 5 formalizes the *partial rewriting* mechanism, which is based on tree simulation. Section 6 introduces our verification technique for detecting incorrect as well as missing/incomplete Web pages. A description of the system GVERDI is provided in Sect. 7. Section 8 concludes. Proofs of the technical results can be found in the appendix.

## 2 Preliminaries

We call a finite set of symbols *as alphabet*. Given the alphabet  $A$ ,  $A^*$  denotes the set of all finite sequences of elements over  $A$ . Syntactic equality between objects is represented by  $\equiv$ .

By  $\mathcal{V}$  we denote a countably infinite set of variables and  $\Sigma$  denotes a set of function symbols, or *signature*. We consider varyadic signatures as in [15] (i.e., signatures in which symbols have an unbounded arity, that is, they may be followed by an arbitrary number of arguments).  $\tau(\Sigma, \mathcal{V})$  and  $\tau(\Sigma)$  denote the *non-ground term algebra* and the *term algebra* built on  $\Sigma \cup \mathcal{V}$  and  $\Sigma$ , respectively.

Terms are viewed as labeled trees in the following (non-standard) way: a term in  $\tau(\Sigma, \mathcal{V})$  is a tree  $(V, E, r, \text{label})$ , where  $V$  is a set of vertices,  $E$  is a set of edges (i.e., pairs of vertices),  $r \in V$  is the *root* vertex and label is a *labeling* function such that  $\text{label}(v) \in (\Sigma \cup \mathcal{V})$ , for each  $v \in V$ . Let us see a small example.

*Example 1* Consider the term  $t \equiv f(h(a), X)$  in  $\tau(\{f, h, a\}, \{X\})$ , which is illustrated in Fig. 1. Term  $t$  can be

represented by the structure  $(V, E, r, \text{label})$ , where  $V \equiv \{v_0, v_1, v_2, v_3\}$ ,  $E \equiv \{(v_0, v_1), (v_0, v_2), (v_1, v_3)\}$ ,  $r \equiv v_0$ , and function label is defined as follows:  $\text{label}(v_0) = f$ ,  $\text{label}(v_1) = h$ ,  $\text{label}(v_2) = X$ ,  $\text{label}(v_3) = a$ .

Given two vertices  $v, v' \in V$  of a term  $t \equiv (V, E, r, \text{label})$ , by  $v \geq v'$  we mean that  $v$  is a descendant of  $v'$  in  $t$ . By  $t|_v$  we mean the subterm rooted at vertex  $v$  of  $t$ .  $t[r]_v$  is the term  $t$  with the subterm rooted at vertex  $v$  replaced by term  $r$ .

We denote the depth of a vertex  $v$  in a term  $t$ , that is the number of edges between  $r$  and  $v$  in  $t$ , as  $\text{depth}(t, v)$ . A substitution  $\sigma \equiv \{X_1/t_1, X_2/t_2, \dots, X_n/t_n\}$  is a mapping from the set of variables  $\mathcal{V}$  into the set of terms  $\tau(\Sigma, \mathcal{V})$  satisfying the following conditions: (1)  $X_i \neq X_j$ , whenever  $i \neq j$ , (2)  $X_i\sigma = t_i$ ,  $i = 1, \dots, n$ , and (3)  $X\sigma = X$ , for any  $X \in \mathcal{V} \setminus \{X_1, \dots, X_n\}$ . By  $\text{Var}(s)$  we denote the set of variables occurring in the syntactic object  $s$ .

In the following, we consider marked terms. Given  $\Sigma$  and  $\mathcal{V}$ , we denote the marked version of  $\Sigma$  ( $\mathcal{V}$ , respectively) as  $\underline{\Sigma}$  ( $\underline{\mathcal{V}}$ , respectively). A syntactic object  $\underline{o} \in \underline{\Sigma} \cup \underline{\mathcal{V}}$  is called the marked version of  $o \in \Sigma \cup \mathcal{V}$ . Given a term  $t \equiv (V, E, r, \text{label}) \in \tau(\Sigma, \mathcal{V})$ , a marking for  $t$  is a (boolean) function  $\mu: V \rightarrow \{\text{yes}, \text{no}\}$ . The empty marking  $\varepsilon$  for  $t$  is a marking for  $t$ , such that  $\varepsilon(v) = \text{no}$ , for each  $v \in V$ . We define the marked part of a term  $t$  as

$$\text{mark}(t, \mu) \equiv (\{v \in V \mid \mu(v) = \text{yes}\}, \{(v_1, v_2) \in E \mid \mu(v_1) = \mu(v_2) = \text{yes}\}, r, \text{label}).$$

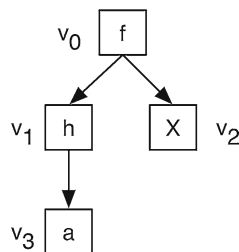
A valid marking  $\mu$  for a term  $t \equiv (V, E, r, \text{label})$  is the empty marking for  $t$  or a marking for  $t$  such that the two following conditions hold:

1.  $\mu(r) = \text{yes}$ .
2.  $\text{mark}(t, \mu)$  is a term in  $\tau(\Sigma, \mathcal{V})$ .

Given a term  $t \equiv (V, E, r, \text{label})$  and a valid marking  $\mu$  for  $t$ , by slightly abusing notation, a marked term  $\mu(t)$  is a term in  $\tau(\underline{\Sigma} \cup \underline{\mathcal{V}}, \underline{\mathcal{V}} \cup \underline{\mathcal{V}})$  such that, for each vertex  $v \in V$ , the label associated with  $v$  in  $t$  is replaced by its marked version in  $\mu(t)$ , whenever  $\mu(v) = \text{yes}$ .

When no confusion can arise, we simply denote the marked term  $\varepsilon(t)$  by  $t$ .

**Fig. 1** Term representation of  $f(h(a), X)$



**Example 2** Consider again the term  $t \equiv (f(h(a), X))$  of Example 1. Let  $\mu_1$  be a marking for  $t$  defined as  $\mu_1(v_0) = \mu_1(v_2) = \mu_1(v_3) = \text{yes}$ ,  $\mu_1(v_1) = \text{no}$ . Additionally, let  $\mu_2$  be a marking for  $t$  such that  $\mu_2(v_0) = \mu_2(v_1) = \text{yes}$ ,  $\mu_2(v_2) = \mu_2(v_3) = \text{no}$ . Note that  $\mu_1$  is not a valid marking for  $t$  as the marked part of  $t$  is not a term in  $\tau(\{f, h, a\}, \{X\})$  [see Fig. 2a], whereas  $\mu_2$  is valid for  $t$  and  $\mu_2(t) = f(h(a), X)$  is a marked term [see Fig. 2b].

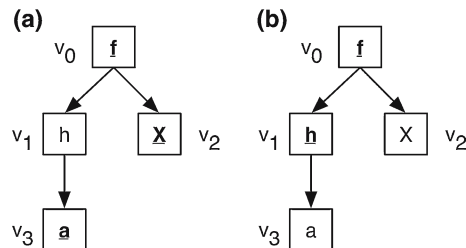
### 2.1 Term rewriting systems

Term rewriting systems provide an adequate computational model for functional languages. In the sequel, we follow the standard framework of term rewriting (see [4,30]). A term rewriting system (TRS for short) is a pair  $(\Sigma, R)$ , where  $\Sigma$  is a signature and  $R$  is a finite set of reduction (or rewrite) rules of the form  $\lambda \rightarrow \rho$ ,  $\lambda, \rho \in \tau(\Sigma, \mathcal{V})$ ,  $\lambda \notin \mathcal{V}$ , and  $\text{Var}(\rho) \subseteq \text{Var}(\lambda)$ . We will often write just  $R$  instead of  $(\Sigma, R)$ . Sometimes, we denote the signature of a TRS  $(\Sigma, R)$  by  $\Sigma_R$ .

We formalize the rewriting relation for  $(V, E, r, \text{label})$  terms as follows. A rewrite step is the application of a rewrite rule to an expression. A term  $s \equiv (V, E, r, \text{label})$  rewrites to a term  $t$  via  $r \in R$ ,  $s \rightarrow_r t$  (or  $s \rightarrow_R t$ ), if there exist  $v \in V$ ,  $r \equiv \lambda \rightarrow \rho$ , and substitution  $\sigma$  such that  $s|_v \equiv \lambda\sigma$  and  $t \equiv s[\rho\sigma]_v$ . When no confusion can arise, we will omit any subscript (i.e.,  $s \rightarrow t$ ). By  $\rightarrow_R^*$ , we denote the reflexive, and transitive closure of  $\rightarrow_R$ . A term  $s$  is a *irreducible form* (or *normal form*) w.r.t.  $R$ , if there is no term  $t$  s.t.  $s \rightarrow_R t$ .  $t$  is the irreducible form of  $s$  w.r.t.  $R$  (in symbols  $s \rightarrow_R^! t$ ) if  $s \rightarrow_R^* t$  and  $t$  is irreducible.

We say that a TRS  $R$  is *terminating*, if there exists no infinite rewrite sequence  $t_1 \rightarrow_R t_2 \rightarrow_R \dots$ . A TRS  $R$  is *confluent* if, for all terms  $s, t_1, t_2$ ,  $s \rightarrow_R^* t_1$  and  $s \rightarrow_R^* t_2$  imply there exists term  $t$  s.t.  $t_1 \rightarrow_R^* t$  and  $t_2 \rightarrow_R^* t$ . When  $R$  is terminating and confluent it is called *canonical*. In canonical TRSs, each input term  $t$  can be univocally reduced to an *irreducible form*.

Let  $s = t$  be an equation, we say that the equation  $s = t$  holds in a canonical TRS  $R$ , if there exists  $z \in \tau(\Sigma, \mathcal{V})$  such that  $s \rightarrow_R^! z$  and  $t \rightarrow_R^! z$ .



**Fig. 2** Examples of valid and non-valid markings for the term  $f(h(a), X)$

**Example 3** Let  $R$  be the following canonical TRS

$$\begin{aligned} \text{sum}(X, 0) &\rightarrow X \\ \text{sum}(X, s(Y)) &\rightarrow s(\text{sum}(X, Y)) \end{aligned}$$

$$\begin{aligned} \text{append}([], L_1) &\rightarrow L_1 \\ \text{append}([X|L_1], L_2) &\rightarrow [X|\text{append}(L_1, L_2)] \end{aligned}$$

$$\begin{aligned} \leq(0, X) &\rightarrow \text{true} \\ \leq(s(X), 0) &\rightarrow \text{false} \\ \leq(s(X), s(Y)) &\rightarrow \leq(X, Y) \end{aligned}$$

In the TRS  $R$  we define three functions. Function *sum* computes the sum of two naturals, function *append* concatenates two lists, and, function  $\leq$  which defines the “less or equal to” relation between natural numbers. Natural numbers are represented in Peano’s notation by means of the constant 0 and the unary operator  $s$ . By abuse we will write  $n \in \mathbb{N}$  as a shorthand for  $s^n(0)$ . We use the standard notation for lists with  $[ ]$  being the empty list. Strings are viewed as lists of characters as usual.

### 3 Denotation of Web sites

In our framework, a *Web page* is either an XML [36] or an XHTML [38] document, which we assumed to be well-formed. There are already programs and online services such as Tidy [33], which are able to validate and correct the XHTML/XML syntax, and Doctor HTML [28], which also performs link-checking.

Let us consider two alphabets  $T$  and  $\text{Tag}$ . We denote the set  $T^*$  by  $\text{Text}$ . An object  $t \in \text{Tag}$  is called *tag* element, while an element  $w \in \text{Text}$  is called *text* element. Since Web pages are provided with a tree-like structure, they can be straightforwardly translated into ordinary terms of a given term algebra  $\tau(\text{Text} \cup \text{Tag})$  as shown in Fig. 3. Note that XML/XHTML tag attributes can be considered as common tagged elements, and hence translated in the same way.

A *marked Web page* is defined as  $\mu(p)$ , where  $p \in \tau(\text{Text} \cup \text{Tag})$  and  $\mu$  is a valid marking for  $p$ . A *Web site* is a finite collection of marked Web pages  $\{\varepsilon(p_1) \dots \varepsilon(p_n)\}$ . In the following, we will also consider terms of the non-ground term algebra  $\tau(\text{Text} \cup \text{Tag}, \mathcal{V})$ , which may contain variables. An element  $s \in \tau(\text{Text} \cup \text{Tag}, \mathcal{V})$  is called *Web page template*.  $\mu(s)$  is a *marked Web page template*, when  $s \in \tau(\text{Text} \cup \text{Tag}, \mathcal{V})$  and  $\mu$  is a valid marking for  $s$ . In our methodology, (marked) Web page templates are used for specifying properties on Web sites as described in the following section.

**Example 4** In the following, we present a Web site  $W$  of a research group, which contains information about group members affiliation, scientific publications, research projects, teaching and personal data.

```
{ (1) members(member(name(mario),
                    surname(rossi),
                    status(professor)),
              member(name(franca),
                    surname(bianchi),
                    status(technician)),
              member(name(giulio),
                    surname(verdi),
                    status(student)),
              member(name(mario),
                    surname(rossi),
                    status(professor))
            ),
  (2) hpage(fullname(mariorossi), phone(3333),
            status(professor), hobbies(
              hobby(reading),
              hobby(gardening))),
  (3) hpage(fullname(francabianchi),
            status(technician), phone(5555),
            links(link(url(www.google.com),
                      urlname(google)),
                  link(url(www.sexycalculus.com),
                      urlname(FormalMethods))),
  (4) hpage(fullname(annagialli),
            status(professor),
            blink(phone(4444)),
            teaching(course(link(
              url(http://www.algebra.math),
              urlname(Algebra)))))
  (5) pubs(pub(name(mario), surname(rossi),
              title(blah1), year(2003)),
          pub(name(anna), surname(gialli),
              title(blah2), year(2002))),
  (6) projects(project(pname(A1), grant1(1000),
                      grant2(200), total(1200), coordinator(
                        fullname(mariorossi))),
              project(pname(B1), grant1(800),
                      grant2(300),
                      projectleader(surname(gialli),
                                    name(anna)),
                      total(1000))) }
```

### 4 Web specification language

In the following, we present a term rewriting specification language, which is helpful to express properties about the content and the structure of a given Web site. Roughly speaking, a Web specification is a pair of finite set of rules. The first set of rules describes constraints for detecting erroneous Web pages (*correctness rules*) as well as discovering incomplete/missing Web pages (*completeness rules*). Diagnoses are carried out

**Fig. 3** An XML document and its corresponding encoding as a ground term  $p$

<pre> &lt;members&gt;   &lt;member status="professor"&gt;     &lt;name&gt; mario &lt;/name&gt;     &lt;surname&gt; rossi &lt;/surname&gt;   &lt;/member&gt;   &lt;member status="technician"&gt;     &lt;name&gt; franca &lt;/name&gt;     &lt;surname&gt; bianchi &lt;/surname&gt;   &lt;/member&gt;   &lt;member status="student"&gt;     &lt;name&gt; giulio &lt;/name&gt;     &lt;surname&gt; verdi &lt;/surname&gt;   &lt;/member&gt; &lt;/members&gt;         </pre>	<pre> members(   member(status(professor),     name(mario),     surname(rossi)   ),   member(status(technician),     name(franca),     surname(bianchi)   )   member(status(student),     name(giulio),     surname(verdi)   ) )         </pre>	
--	---	--

by running Web specifications on Web sites. The operational mechanism, formalized in Sect. 5, is based on a novel rewriting-based technique, which is able to extract partial structure from a term, and then rewrite it. The second set of rules  $R$  contains the definition of some auxiliary functions which the user would like to provide, such as string processing, arithmetic, boolean operators, etc. It is formalized as a canonical term rewriting system which is handled by standard rewriting [30].

Correctness rules allow us to recognize incorrect and forbidden patterns in the Web site and to locate the wrong Web pages containing these errors. We address these problem by considering the specific nature of semi-structured documents, that is, we combine a structured pattern search technique with a more standard textual search, which is based on regular expressions detection. For the sake of simplicity, we consider an intuitive Unix-like regular expression syntax [34]. In contrast to other Web verification methods, we are able to deal with correctness rules containing interpreted functions, which allow us to check whether the values that may appear in the semistructured documents are correctly computed.

The verification process is carried out by first extracting a partial structure of a (possibly incorrect) Web page, and then checking whether the required constraints are fulfilled. These constraints are expressed by means of equations and regular expressions. Correctness rules are formalized as follows.

**Definition 1** (Correctness rule) *Let  $(\Sigma, R)$  be a canonical TRS. A correctness rule has the following form*

$$l \rightarrow \text{error} \mid C$$

where

1.  $l \in \tau(\text{Text} \cup \text{Tag}, \mathcal{V})$  is a Web page template and  $\text{error} \notin (\text{Text} \cup \text{Tag} \cup \Sigma)$  is a new fresh constant.

2.  $C$  is a (possibly empty) sequence

$$X_1 \text{ in } \text{rexp}_1, \dots, X_n \text{ in } \text{rexp}_n, \Gamma$$

with  $\text{Var}(C) \subseteq \text{Var}(l)$ ,  $\text{rexp}_i$  a regular expression over  $(\text{Text} \cup \text{Tag})$ ,  $i = 1, \dots, n$ , and  $\Gamma$  a sequence of equations over  $\tau(\Sigma, \mathcal{V})$ .

When  $C$  is empty, we simply write  $l \rightarrow \text{error}$ .

Given a correctness rule  $r \equiv (l \rightarrow \text{error} \mid C)$ , we call  $l \rightarrow \text{error}$  the *unconditional part of  $r$*  and we denote it by  $r_u$ . For the sake of expressiveness, we also allow to write inequalities of the form  $s \neq t$  in the conditional part of the correctness rules. Such inequalities are just syntactic sugar for  $(s = t) = \text{false}$ . Informally, the meaning of a correctness rule  $l \rightarrow \text{error} \mid C$  is the following. Whenever an instance  $l\sigma$  of  $l$  is recognized in some Web page  $p$ , and

- (1) Each structured text  $X_i\sigma$ ,  $i = 1, \dots, n$ , is contained in the language of the corresponding regular expression  $\text{rexp}_i$ .
- (2) Each instantiated equality  $(s = t)\sigma$  in  $\Gamma$  holds in the canonical TRS  $R$ .

then, Web page  $p$  is signaled as an incorrect page.

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  $\langle A \rangle$  and  $\langle E \rangle$  to distinguish “universal” from “existential” rules. Right-hand sides of completeness rules may contain functions, which are defined by the user via a canonical TRS.

**Definition 2** (Completeness rule) *Let  $(\Sigma, R)$  be a canonical TRS. A completeness rule is either a universal rule of the form  $l \rightarrow \mu(x) \langle A \rangle$  or an existential rule of the form  $l \rightarrow \mu(x) \langle E \rangle$ , where  $l \in \tau(\text{Text} \cup \text{Tag}, \mathcal{V})$ ,*

$x \in \tau(\text{Text} \cup \text{Tag} \cup \Sigma, \mathcal{V})$ ,  $\text{Var}(x) \subseteq \text{Var}(1)$  and  $\mu$  is a valid marking for  $x$ .

Intuitively, the interpretation of a universal rule  $1 \rightarrow \mu(x) \langle A \rangle$  [respectively, an existential rule  $1 \rightarrow \mu(x) \langle E \rangle$ ] w.r.t. a Web site  $\mathbb{W}$  is as follows: if (an instance of)  $1$  is recognized in  $\mathbb{W}$ , also (an instance of) the irreducible form of  $x$  must be recognized in *all* (respectively, *some*) of the Web pages which embed (an instance of) the marked part of  $x$ . Informally, marking information provides the “scope” of the universal/existential quantifiers of the rule, and allows us to compute the part of the Web site which is affected by the quantification.

*Example 5* Consider the Web site

$$\mathbb{W} = \{f(h(g(a), t)), h(t, s(m, n), g(p)), h(n, p(r, q)), w(s(m, n))\}$$

and the rule  $r \equiv (w(x) \rightarrow \underline{h}(t, x) \langle q \rangle)$ ,  $q \in \{E, A\}$ . Then, the marked part of the right-hand side of  $r$  is  $h(t)$ .

Thus, the property, which is modeled by  $r$ , will be checked only on the set  $\{f(h(g(a), t)), h(t, s(m, n), g(p))\}$ , which contains all and only the Web pages of  $\mathbb{W}$  embedding the term  $h(t)$ .

Formally, Web specifications are as follows.

**Definition 3** (Web specification) *A Web specification is a pair  $(\mathbb{I}, R)$ , where  $\mathbb{I} \equiv \mathbb{I}_N \uplus \mathbb{I}_M$  is a finite set of rules such that  $\mathbb{I}_N$  (respectively,  $\mathbb{I}_M$ ) is a set of correctness (respectively, completeness) rules, and  $R$  is a canonical TRS.*

In order to provide an effective, and terminating verification method, we require that the arguments of the function calls in the right-hand sides of the completeness rules do not contain symbols in

$$\{f \mid f(t_1, \dots, t_n) \rightarrow \mu(x) \in \mathbb{I}_M\}$$

(for the technical reasons, see Lemma 3 in the appendix).

Given a set of completeness rules  $\mathbb{I}_M$ , we denote the set of all left-hand sides (right-hand sides without marks, respectively) of rules in  $\mathbb{I}_M$  by  $\text{Lhs}_M$  ( $\text{Rhs}_M$ , respectively). In symbols,  $\text{Lhs}_M = \{1 \mid 1 \rightarrow \mu(x) \in \mathbb{I}_M\}$  and  $\text{Rhs}_M = \{x \mid 1 \rightarrow \mu(x) \in \mathbb{I}_M\}$ .

Note that, by using the traditional encoding of boolean operations by means of rewrite rules [31, 23], and by introducing non-deterministic rewriting [27], it would also be possible to extend our framework to a richer specification language providing non-confluent functions such as those that express disjunctive conditions (so that the presence of information of kind  $A$  requires information of kind  $B$  or  $C$  to be also represented). For the sake of simplicity, we do not deal with non-confluent TRSs in this work.

The following example illustrates the definition of a Web specification. Marks are introduced by the user to select those Web pages for which she wants to check the specified integrity conditions.

*Example 6* Consider the Web specification which consists of the canonical TRS  $R$  of Example 3 together with the built-in definition of function  $\text{Nat}(X)$ , which translates a string  $x$  to a natural number, and the completeness and correctness rules of Fig. 4.

This Web specification models some required properties for the Web site of Example 4.

The first four rules are completeness rules, while the remaining ones are correctness rules. First rule formalizes the following property: if there is a Web page containing a member list, then for each member, a home page should exist which contains (at least) the full name and the status of this member. The full name is computed by appending the name and the surname strings by means of the standard `append` function whose definition is given in TRS  $R$ . The marking information establishes that the property must be checked only on home pages (i.e., pages containing the tag “hpage”). Second rule states that, whenever a home page of a professor is recognized, that page must also include some teaching information. The rule is universal, since it must hold for each professor home page. Such home pages are selected by exploiting the mark given on the tag “status”. Third rule specifies that, whenever there exists a Web page containing information about scientific publications, each author of a publication should be a member of the research group. In this case, we must check the property only in the Web page containing the group member list. The fourth rule formalizes that, for each link to a course, a page describing that course must exist. The verification process is carried out only on Web pages containing course information as described by marks.

The fifth rule forbids sexual contents from being published in the home pages of the group members. This is enforced by requiring that the word `sex` does not occur in any home page by using the regular expression  $[:\text{TextTag}:]^* \text{sex} [:\text{TextTag}:]^*$ , which identifies the regular language of all the strings built over  $(\text{Text} \cup \text{Tag})$  containing that word. The sixth rule is provided with the aim of improving accessibility for people with disabilities. It simply states that blinking text is forbidden in the whole Web site. The seventh rule states that, for each research project, the total project budget must be equal to the sum of the funds, which has been granted for the first and the second research period. The eighth rule formalizes the condition that only recent publications are referred to in the Web site (after the

**Fig. 4** Correctness and completeness rules of Example 6

```

member(name(X), surname(Y)) → hpage(fullname(append(X, Y)), status) ⟨E⟩
hpage(status(professor)) → hpage(status(professor), teaching) ⟨A⟩
pubs(pub(name(X), surname(Y))) → member(name(X), surname(Y)) ⟨E⟩
courselink(url(X), urlname(Y)) → cpage(title(Y)) ⟨E⟩
hpage(X) → error | X in [[:TextTag:]]* sex [[:TextTag:]]*
blink(X) → error
project(grant1(X), grant2(Y), total(Z)) → error | sum(Nat(X), Nat(Y)) ≠ Nat(Z)
pub(year(X)) → error | X in [0-9]*, ≤ (Nat(X), 1999) = true
members(member(name(X), surname(Y)), member(name(X), surname(Y))) → error
    
```

year 1999). Finally, the last rule forbids repetitions of the same member entry in the group member list.

The verification of both kinds of rules is mechanized by means of partial rewriting.

### 5 Partial rewriting

In order to mechanize the intended semantics of Web specification rules, we first devise a mechanism which is able to recognize the structure and the labeling of a given Web page template inside a particular page of the Web site. This is provided by page simulation. Without loss of generality, in the following, we disregard the conditional part of correctness rules and the existential/universal quantification of completeness rules. In other words, a rule is just a pair of (possibly) marked terms  $l \rightarrow \mu(x)$ .

#### 5.1 Page simulations

The notion of *page simulation* for Web pages allows us to analyze and extract the partial structure of the Web site which is subject to verification.

Roughly speaking, a Web page  $p_1$  is simulated by a Web page  $p_2$ , if the tree-structure of  $p_1$  is “embedded” into the tree-structure of  $p_2$ . In other words, a simulation of a Web page (i.e., a labeled tree)  $p_1$  in a Web page  $p_2$  can be seen as a relation among the nodes of  $p_1$  and the nodes of  $p_2$  which preserves the edges and the labelings. Before formalizing the idea, we illustrate it by means of a rather intuitive example.

*Example 7* Consider the following Web pages (called  $p_1$  and  $p_2$ , respectively):

```

hpage(name, surname, status(professor), teaching)
hpage(name(mario), surname(rossi),
      status(professor),
      teaching(course(logic1),
                course(logic2)),
      hobbies(hobby(reading), hobby(gardening)))
    
```

Looking at Fig. 5, we observe that the structure of  $p_1$  can be recognized inside the structure of  $p_2$  by considering the relation among nodes of  $p_1$  and  $p_2$ , which is

described by the dashed arrows in the figure. This relation essentially provides the so-called *simulation of  $p_1$  in  $p_2$* . Note that vice-versa does not hold: no relations can be found among nodes of  $p_2$  and  $p_1$ , which “embed” the structure of  $p_2$  into  $p_1$ . In other words, there does not exist a simulation of  $p_2$  in  $p_1$ .

Simulations have been used in a number of works dealing with querying and transformation of semistructured data. For instance, [1, 22] propose some techniques based on simulation for analyzing semistructured data w.r.t. a given schema. The language Xcerpt [11, 10] is a (logic) query language for XML and semistructured documents, which implements a sort of unification by exploiting the notion of graph simulation. Other approaches involving simulation, or closely related notions, have been employed to measure similarity among semistructured documents [8]. To keep our framework simple, we do not consider a semantic change/load for labels; this would require to introduce ontologies, which are outside the scope of the work.

In the following, we provide our notion of simulation which is a slight adaptation of the one given in [10] to consider Web page templates: we generalize the usual label relation to cope with the case when variables are used as labels, in the following definition.

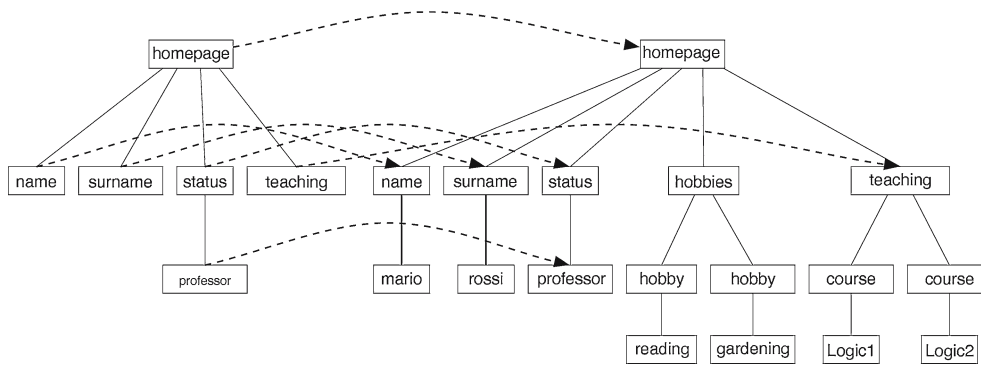
**Definition 4** Let  $s_1 \equiv (V_1, E_1, r_1, label_1)$ ,  $s_2 \equiv (V_2, E_2, r_2, label_2)$  be two Web page templates in  $\tau(\text{Text} \cup \text{Tag}, \mathcal{V})$ . The label relation  $\sim \subseteq V_1 \times V_2$  is defined as follows:

$v_1 \sim v_2$  iff  $label_1(v_1) = label_2(v_2)$  or  $label_1(v_1) \in \mathcal{V}$ .

**Definition 5** Let  $s_1 \equiv (V_1, E_1, r_1, label_1)$ ,  $s_2 \equiv (V_2, E_2, r_2, label_2)$  be two Web page templates in  $\tau(\text{Text} \cup \text{Tag}, \mathcal{V})$  and  $\sim \subseteq V_1 \times V_2$  be the corresponding label relation. A (page) simulation of  $s_1$  in  $s_2$  w.r.t.  $\sim$  is a relation  $\mathbf{S} \subseteq V_1 \times V_2$  such that, for each  $v_1 \in V_1, v_2 \in V_2$

1.  $r_1 \mathbf{S} r_2$ .
2.  $v_1 \mathbf{S} v_2 \Rightarrow v_1 \sim v_2$ .
3.  $v_1 \mathbf{S} v_2 \wedge (v_1, v'_1) \in E_1 \Rightarrow \exists v'_2 \in V_2, v'_1 \mathbf{S} v'_2 \wedge (v_2, v'_2) \in E_2$ .

We define the *projection* of a simulation  $\mathbf{S}$  of  $s_1$  in  $s_2$  w.r.t.  $\sim$  as  $\pi(\mathbf{S}) = \{v_2 \mid (v_1, v_2) \in \mathbf{S}\}$ .



**Fig. 5** Page simulation between  $p_1$  and  $p_2$

Roughly speaking, Definition 5 ensures two degrees of similarity between Web page templates, not only w.r.t. the labelings but also w.r.t. the structures of the templates. On the one hand, Condition (2) of Definition 5 formalizes the similarity w.r.t. labelings, that is, any pair of nodes  $(v, v')$  in a page simulation  $\mathbf{S}$  of  $s_1$  in  $s_2$  have the same label, otherwise node  $v$  must be labeled by a variable, which somehow means that the label of  $v$  can be seen as a generalization of any concrete label of  $v'$ . Finally, Conditions (1) and (3) provide a relation between the tree structure of  $s_1$  and  $s_2$ .

Note that simulations are just relations among nodes of two given Web page templates. For our purposes, we are interested in simulations which are injective mappings from nodes of a given Web page template to nodes of another Web page template. As it will be apparent later, those simulations allow us to project the structure of a Web page template into another one, thus performing a sort of “partial” pattern matching between templates, which will be exploited to formulate our verification technique. So, we first define a subclass of simulations called minimal simulations.

**Definition 6** Let  $s_1 \equiv (V_1, E_1, r_1, \text{label}_1)$ ,  $s_2 \equiv (V_2, E_2, r_2, \text{label}_2)$  be two Web page templates in  $\tau(\text{Text} \cup \text{Tag}, \mathcal{V})$ . A page simulation  $\mathbf{S}$  of  $s_1$  in  $s_2$  w.r.t.  $\sim$  is minimal if there are no page simulations  $\mathbf{S}'$  of  $s_1$  in  $s_2$  w.r.t.  $\sim$  such that  $\mathbf{S}' \subseteq \mathbf{S}$ .

It is straightforward to prove that minimal simulations are mappings.

**Proposition 1** Let  $s_1 \equiv (V_1, E_1, r_1, \text{label}_1)$ ,  $s_2 \equiv (V_2, E_2, r_2, \text{label}_2)$  be two Web page templates in  $\tau(\text{Text} \cup \text{Tag}, \mathcal{V})$ . A minimal page simulation  $\mathbf{S}$  of  $s_1$  in  $s_2$  w.r.t.  $\sim$  is a mapping  $\mathbf{S} : V_1 \rightarrow V_2$ .

Let us see an example which illustrates the notion of minimal simulation.

*Example 8* Let us consider the following Web page templates  $s_1$  and  $s_2$ :

$\text{hobbies}(\text{hobby}(X))$   
 $\text{hobbies}(\text{hobby}(\text{reading}), \text{hobby}(\text{gardening}))$ .

In Fig. 6a, the dashed arrows represent a non-minimal simulation of  $s_1$  in  $s_2$ , while in Fig. 6b and c two minimal simulations of  $s_1$  in  $s_2$  are depicted. Note that the last two simulations are mappings.

However, minimal simulations do not guarantee that the tree structure of a given Web page template is recognized inside another template. Consider, for instance, the page simulation of  $f(X, Y)$  in  $f(a)$  depicted in Fig. 7: it is minimal, but the tree structures of  $f(X, Y)$  and  $f(a)$  are distinct.

For this purpose, we need a one-to-one correspondence between edges of considered Web page templates. Therefore, we only consider minimal and *injective* simulations. Given two Web page templates  $s$  and  $t$ , we denote by  $s \cong t$ , the fact that there exists a minimal, injective simulation of  $s$  in  $t$  w.r.t.  $\sim$ .

It is not difficult to prove that minimal injective simulations are particular instances of Kruskal’s *embeddings* [9] w.r.t. the relation  $\sim$ . In other words, a minimal injective page simulation of  $s_1$  in  $s_2$  w.r.t.  $\sim$  exists iff  $s_1$  is embedded into  $s_2$  w.r.t.  $\sim$ , i.e., we are able to find out the structure and the labeling of  $s_1$  inside  $s_2$ .

### 5.2 Rewriting Web page templates

**Definition 7** Let  $s_1 \equiv (V_1, E_1, r_1, \text{label}_1)$ ,  $s_2 \equiv (V_2, E_2, r_2, \text{label}_2) \in \tau(\text{Text} \cup \text{Tag}, \mathcal{V})$ . We say that  $s_2$  partially matches  $s_1$  via substitution  $\sigma$  iff

1. There exists a minimal, injective page simulation  $\mathbf{S}$  of  $s_1$  in  $s_2$  w.r.t.  $\sim$ .
2. For each  $(v, v') \in \mathbf{S}$  such that  $\text{label}(v) = X \in \mathcal{V}$ ,  $\sigma(X) = (s_2|_{v'})$ .

In Definition 7, we consider only minimal, injective simulations between Web page templates  $s_1$  and  $s_2$  to easily compute a substitution  $\sigma$  such that there exists a



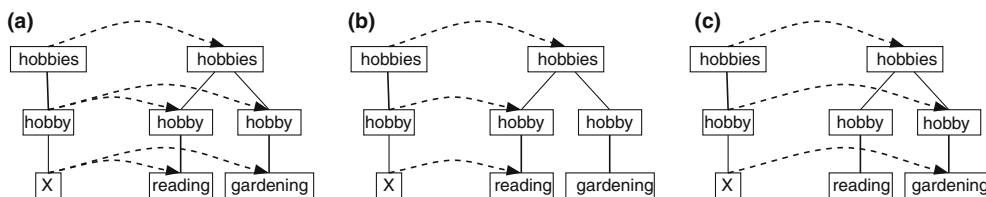


Fig. 6 Non-minimal and minimal simulations

simulation of  $s_1\sigma$  in  $s_2$  w.r.t.  $\sim$ ; in other words,  $s_1\sigma$  is embedded into  $s_2$ . It is worth noting that non-minimal simulations not always ensure the existence of such a substitution. This is the reason why the minimal simulations are also required to be injective. Let us see an example.

*Example 9* Consider again Example 8. We have that  $s_2$  partially matches  $s_1$  via  $\{X/\text{reading}\}$  (see Fig. 6b) and  $s_2$  partially matches  $s_1$  via  $\{X/\text{gardening}\}$  (see Fig. 6c). Note that performing partial matching by the non-minimal simulation of Fig. 6 a would produce  $\sigma \equiv \{X/\text{reading}, X/\text{gardening}\}$ , which is not a substitution.

Now we are ready to define a partial rewrite relation between marked Web page templates, which includes a simplification stage using the user functions of a TRS  $R$ .

**Definition 8** Let  $R$  be a canonical TRS. Let  $s \equiv (V, E, r_s, \text{label})$ ,  $t \in \tau(\text{Text} \cup \text{Tag}, \mathcal{V})$ . Let  $\mu_1$  and  $\mu_2$  be two valid markings for  $s$  and  $t$ , respectively. Then,  $\mu_1(s)$  partially rewrites to  $\mu_2(t)$  via rule  $rl \equiv l \rightarrow \mu(r)$  and substitution  $\sigma$  (in symbols,  $\mu_1(s) \rightarrow_{rl}^{\sigma} \mu_2(t)$ ) iff there exists  $v \in V$  such that

1.  $s|_v$  partially matches  $l$  via  $\sigma$ .
2. Let  $r \equiv (V_r, E_r, r, \text{label}_r)$  and  $r\sigma \equiv (V_{r\sigma}, E_{r\sigma}, r, \text{label}_{r\sigma})$ . For each  $v \in V_{r\sigma}$ ,

$$\mu_2(v) = \begin{cases} \mu(v) & \text{if } v \in (V_r \cap V_{r\sigma}) \\ \mu(v') & \text{if } v \in (V_{r\sigma} \setminus V_r) \wedge (\exists v' \in V_r, v \geq v', \text{label}_r(v') \in \text{Var}(r)) \end{cases}$$

3.  $t = \text{Reduce}(r\sigma, R)$ , where function  $\text{Reduce}(x, R)$  computes, by standard term rewriting, the irreducible form of  $x$  in  $R$  ignoring the eventual marks for the functions in  $R$ .

When rule  $rl$  and substitution  $\sigma$  are understood, we simply write  $\mu_1(s) \rightarrow \mu_2(t)$ .

It is worth noting that we provide a notion of partial rewriting in which the context of the selected reducible expression  $s|_v$  of the Web page template which is rewritten is disregarded after the rewrite step [see point (3) of

Definition 8]. Roughly speaking, given a Web specification rule  $l \rightarrow \mu(r)$ , partial rewriting allows us to extract a subpart of a given Web page (template)  $s$ , which partially matches  $l$ , and to replace  $s$  by a reduced instance of  $r$ ; namely,  $\text{Reduce}(r\sigma, R)$  [see points (1) and (3) of Definition 8]. Point (2) of Definition 8 establishes that rewritten templates inherit marks from the right-hand sides of the applied rules. More precisely,

- Each vertex of  $r\sigma$ , which is not affected by substitution  $\sigma$ , maintains the same marking of  $r$ .
- Each vertex, which belongs to a subterm of  $r\sigma$  replacing a variable  $X$  of  $r$ , is marked yes.
- Each vertex, which belongs to a subterm of  $r\sigma$  replacing a variable  $X$  of  $r$ , is marked no.

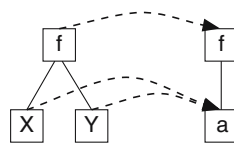
A partial rewrite sequence is of the form  $\mu_0(s_0) \rightarrow_{r_0}^{\sigma_0} \mu_1(s_1) \rightarrow_{r_1}^{\sigma_1} \dots$ . Moreover, we denote the transitive closure (resp., the transitive and reflexive closure) of  $\rightarrow$  by  $\rightarrow^+$  (resp.,  $\rightarrow^*$ ). By notation  $\mu_0(s_0) \rightarrow^n \mu_1(s_1)$  we denote a partial rewrite sequence of length  $n$  (i.e., a partial rewrite sequence which is made up of  $n$  partial rewrite steps).

*Example 10* Consider the Web page  $p$  of Fig. 3 and the first rule  $r_1$  of the Web specification of Example 6. Let us suppose that TRS  $R$  defines the standard function  $\text{append}$  for concatenating strings. Then, Web page template  $\varepsilon(p)$  partially rewrites to the following Web pages by applying  $r_1$ .

$$\begin{aligned} \varepsilon(p) \rightarrow_{r_1} \text{Reduce}(\text{hpage}(\text{fullname}(\text{append}(\text{mario}, \text{rossi})), \text{status}), R) &= \text{hpage}(\text{fullname}(\text{mariorossi}), \text{status}) \\ \varepsilon(p) \rightarrow_{r_1} \text{Reduce}(\text{hpage}(\text{fullname}(\text{append}(\text{franca}, \text{bianchi})), \text{status}), R) &= \text{hpage}(\text{fullname}(\text{francabianchi}), \text{status}) \\ \varepsilon(p) \rightarrow_{r_1} \text{Reduce}(\text{hpage}(\text{fullname}(\text{append}(\text{giulio}, \text{verdi})), \text{status}), R) &= \text{hpage}(\text{fullname}(\text{giulioverdi}), \text{status}) \end{aligned}$$

Roughly speaking, marks in the right-hand sides of the rules allow us to find sets of Web pages, which might

**Fig. 7** Minimal, non-injective simulation



be incomplete or missing. Then, real buggy pages are detected inside these sets. We formalize the idea in the following section.

**6 The verification framework**

In the following, we show how simulation and partial rewriting can be applied to verify a given Web site w.r.t. a Web specification. As we have seen in Sect. 5.1, simulation allows us to identify the structure of a given Web page (possibly, a template) into another one. By taking advantage of this fact, we can develop a methodology, which is able to discover correctness as well as completeness errors in a given Web site w.r.t. a Web specification.

More precisely, our analysis allows us to discover the following kinds of errors:

- Erroneous/forbidden information in the Web site (*correctness errors*).
- Web pages which are missing in a Web site or Web pages which are incomplete w.r.t. a given Web specification (*completeness errors*).

**6.1 Detecting correctness errors**

In this section, we provide a simple mechanism based on partial rewriting which can detect erroneous or undesirable data included in a Web site. Our methodology allows us to precisely locate which part of a Web page does not fulfill the Web specification. We apply correctness rules to the Web pages of the Web site in order to discover incorrect patterns. More precisely, given a Web page  $p$ , we first try to recognize a given Web page template  $l$  into  $p$  by partially rewriting  $p$  via the (unconditional part of a) correctness rule  $l \rightarrow error | C$ . Then, we analyze the values taken by the variables of  $C$ , which are obtained as a by-product of the partial rewrite step. If the structured text, which is bound to each variable in  $C$ , belongs to the language of the corresponding regular expression, and all the instantiated equations in  $C$  hold, the faulty Web page  $p$  and an incorrectness symptom are supplied to the user.

**Definition 9** Let  $W$  be a Web site,  $(I_N \uplus I_M, R)$  be a Web specification. Given  $p \in W$ , we say that  $p$  is incorrect w.r.t.  $(I_N \uplus I_M, R)$ , if there exists a rule  $r \equiv (l \rightarrow error | X_1$  in  $rexp_1, \dots, X_n$  in  $rexp_n, \Gamma) \in I_N$  such that

1.  $p \rightarrow_{r_u}^\sigma error$ , where  $r_u$  is the unconditional part of  $r$ .
2. For  $i = 1, \dots, n$ ,  $X_i \sigma \in \mathcal{L}(rexp_i)$ , where  $\mathcal{L}(rexp_i)$  is the regular language described by  $rexp_i$ .
3. Each instantiated equation of  $\Gamma$ ,  $(s = t)\sigma$ , holds in  $R$ .

We say that  $l\sigma$  is an incorrectness symptom for  $p$ .

Let us see an example which illustrates the above definition.

*Example 11* Let  $(I_N \uplus I_M, R)$  be the Web specification of Example 6 and  $W$  be the Web site of Example 4.

Now, consider the correctness rule

$$r \equiv hpage(X) \rightarrow error | X \text{ in } [:TextTag:] * sex[:TextTag:] * .$$

Note that the only Web page in  $W$  which can yield a correctness error by using  $r$  is Web page (3), since (3) can be partially rewritten to  $error$  via  $r_u$  by means of the following substitution  $\sigma$

$$\{X/links(link(url(www.google.com), urlname(google)), link(url(www.sexycalculus.com), urlname(FormalMethods)))\}$$

and  $X\sigma$  belongs to

$$\mathcal{L}([:TextTag:] * sex[:TextTag:] *).$$

The corresponding incorrectness symptom is

$$hp(links(link(url(www.google.com), urlname(google)), link(url(www.sexycalculus.com), urlname(FormalMethods)))) .$$

Web page (4) is incorrect w.r.t.  $(I_N \uplus I_M, R)$ , as it rewrites to  $error$  by rule  $blink(X) \rightarrow error$  and gives rise to the incorrectness symptom  $blink(phone(4444))$ . Moreover, we can also discover that the total budget of project B1 is wrongly computed by applying rule

$$project(grant1(X), grant2(Y), total(Z)) \rightarrow error | sum(Nat(X), Nat(Y)) \neq Nat(Z).$$

Indeed, Web page (6) is incorrect, since the equation

$$sum(Nat(800), Nat(300)) \neq Nat(1000)$$

hold in  $R$ . The associated incorrectness symptom is

$$project(grant1(800), grant2(300), total(1000)).$$

Finally, Rule

$$r' \equiv \text{members}(\text{member}(\text{name}(X), \text{surname}(Y)), \text{member}(\text{name}(X), \text{surname}(Y))) \rightarrow \text{error}$$

detects that the entry for Mario Rossi in Web page (1) is repeated twice, since it derives error by using  $r'$  via substitution  $\{X/\text{mario}, Y/\text{rossi}\}$ . The respective incorrectness symptom is:

$$\text{members}(\text{member}(\text{name}(\text{mario}), \text{surname}(\text{rossi})), \text{member}(\text{name}(\text{mario}), \text{surname}(\text{rossi}))).$$

The example above points out the usefulness of incorrectness symptoms: they allow us to precisely locate which erroneous piece of information must be modified by the user in order to repair the faulty Web site.

Algorithm 1 outlines a procedure for the detection of correctness errors, which takes as input a Web site  $\mathbb{W}$ , a set of correctness rules  $\mathbb{I}_N$ , and a canonical TRS  $R$ . Basically, the procedure repeatedly applies the test of Definition 9 for recognizing incorrect and forbidden patterns. More precisely, for every Web page in  $\mathbb{W}$ , we verify whether (1)  $p$  reduces to the constant error via the unconditional part of some correctness rule  $r$ , and (2) the constraints in the condition of  $r$  are fulfilled. In the case that an error is found in a Web page  $p$  by using a rule  $l \rightarrow r \mid C$ , the pair  $(p, l\sigma)$  is returned, which consists of the wrong page together with the corresponding incorrectness symptom.

**Algorithm 1** An algorithm for detecting correctness errors in a Web site.

```

1: procedure CORRECTNESS-ERRORS( $\mathbb{W}, \mathbb{I}_N, R$ )
2:   for all  $p \in \mathbb{W}$  do
3:     for all  $r \equiv (l \rightarrow \text{error} \mid$ 
          $X_1 \text{ in } \text{rexp}_1, \dots, X_n \text{ in } \text{rexp}_n,$ 
          $s_1 = t_1, \dots, s_m = t_m) \in \mathbb{I}_N$  do
4:       if  $(p \rightarrow_r^\sigma \text{error})$  then
5:         if  $(X_i \sigma \in \mathcal{L}(\text{rexp}_i), i = 1, \dots, n)$  and
            $((s_j = t_j)\sigma, j = 1, \dots, m, \text{ holds in } R)$  then
6:           output ("Error:  $(p, l\sigma)$ ") end if
7:         end if
8:       end for
9:     end for
10: end procedure

```

**Proposition 2** Let  $\mathbb{W}$  be a Web site, and  $(\mathbb{I}_N \uplus \mathbb{I}_M, R)$  be a Web specification. Then, the procedure CORRECTNESS-ERRORS( $\mathbb{W}, \mathbb{I}_N, R$ ) terminates, and for each pair  $(p, l\sigma)$ , which is returned,  $p$  is an incorrect Web page w.r.t.  $(\mathbb{I}_N \uplus \mathbb{I}_M, R)$  and  $l\sigma$  is the corresponding incorrectness symptom.

### 6.2 Detecting completeness errors

Essentially, the main idea to diagnose completeness errors is to compute the set of all possible marked

expressions that can be derived from  $\mathbb{W}$  via the completeness rules of a Web specification  $(\mathbb{I}_N \uplus \mathbb{I}_M, R)$  by means of partial rewriting. These marked terms can be thought of as requirements to be fulfilled by  $\mathbb{W}$ . Then, we check whether the computed requirements are satisfied by  $\mathbb{W}$  using simulation and marking/quantification information.

In summary, the method works in two steps, as described below.

1. Compute the set of requirements  $\text{Req}_{\mathbb{M}, \mathbb{W}}$  for  $\mathbb{W}$  w.r.t.  $\mathbb{I}_M$ .
2. Check  $\text{Req}_{\mathbb{M}, \mathbb{W}}$  in  $\mathbb{W}$ .

Formally, a *requirement* is a pair  $\langle \mu(e), \mathfrak{q} \rangle$ , where  $\mu(e)$  is a marked term and  $\mathfrak{q} \in \{A, E, \_ \}$ . A requirement is called *universal* whenever  $\mathfrak{q} = A$ , while it is called *existential* when  $\mathfrak{q} = E$ . Requirements of the form  $\langle \mu(e), \_ \rangle$  are called *non-quantified* requirements.

In order to formalize step 1, we define the following operator.

**Definition 10** Let  $\mathbb{T}$  be a set of marked terms and  $(\mathbb{I}_N \uplus \mathbb{I}_M, R)$  be a Web specification. Then, the immediate completeness requirements operator

$$\mathcal{J}_M(\bar{\mathbb{T}}) = \bar{\mathbb{T}} \cup \{ \langle \mu_2(s_2), \mathfrak{q} \rangle \mid \exists \langle \mu_1(s_1), \mathfrak{q}_1 \rangle \in \bar{\mathbb{T}}, r \equiv (l \rightarrow \mu(r) \langle \mathfrak{q} \rangle) \in \mathbb{I}_M \text{ s.t. } \mu_1(s_1) \rightarrow_r \mu_2(s_2) \}$$

where  $\bar{\mathbb{T}} = \{ \langle s, \_ \rangle \mid s \in \mathbb{T} \}$ .

Roughly speaking, the operator in Definition 10 computes all the requirements which are obtained by partially rewriting the marked expressions in  $\bar{\mathbb{T}}$  using the completeness rules of  $\mathbb{I}_M$ , and returns the union of the resulting set and  $\bar{\mathbb{T}}$ . By repeatedly applying this operator, it is possible to compute all marked terms that can be derived from an initial Web site after an arbitrary number of partial rewriting steps. For this purpose, we formalize the *ordinal powers* of the operator  $\mathcal{J}_M$  w.r.t. a Web site  $\mathbb{W}$  as follows:  $\mathcal{J}_M \uparrow^{\mathbb{W}} 0 = \bar{\mathbb{W}}, \mathcal{J}_M \uparrow^{\mathbb{W}} n = \mathcal{J}_M(\mathcal{J}_M \uparrow^{\mathbb{W}} (n - 1)), n > 0$ .

It is immediate to prove that the operator  $\mathcal{J}_M$  is continuous on the lattice consisting of the powerset of the requirements ordered by set inclusion. This ensures that a least fixpoint of  $\mathcal{J}_M$  exists and can be reached after  $\omega$  applications of  $\mathcal{J}_M$ , that is,  $\mathcal{J}_M \uparrow^{\mathbb{W}} \omega$  where  $\omega$  is the first infinite ordinal. The least fixpoint of  $\mathcal{J}_M$  contains all the marked expressions derivable from  $\mathbb{W}$  via  $\mathbb{I}_M$  along with their quantification information.

Now, recalling the interpretation of the completeness rules of the Web site specification given in Sect. 4, marked terms derived by the application of a completeness rule must be recognized as (part of) some Web page

in the Web site. Therefore, the expressions in the least fixpoint of  $\mathcal{J}_M$  provide information that must occur in  $W$ . Thus, since the pages in  $W$  trivially occur in the Web site  $W$ , we define the *set of requirements* for  $W$  w.r.t.  $I_M$  as

$$Req_{M,W} = lfp(\mathcal{J}_M) \setminus \bar{W}$$

where  $lfp(\mathcal{J}_M)$  is the least fixpoint of the operator  $\mathcal{J}_M$ . Observe that the set  $Req_{M,W}$  does not contain any non-quantified requirement.

*Example 12* Consider the Web specification  $(I_N \uplus I_M, R)$  of Example 6 and the Web site  $W$  of Example 4. Then, the set of computed requirements  $Req_{M,W}$  is

$$\{ \langle \underline{hpage}(\underline{fullname}(\underline{mario} \underline{rossi}), \underline{status}), E \rangle, \langle \underline{hpage}(\underline{fullname}(\underline{francabianchi}), \underline{status}), E \rangle, \langle \underline{hpage}(\underline{fullname}(\underline{giulio} \underline{verdi}), \underline{status}), E \rangle, \langle \underline{hpage}(\underline{status}(\underline{professor}), \underline{teaching}), A \rangle, \langle \underline{member}(\underline{name}(\underline{mario}), \underline{surname}(\underline{rossi})), E \rangle, \langle \underline{member}(\underline{name}(\underline{anna}), \underline{surname}(\underline{gialli})), E \rangle, \langle \underline{hpage}(\underline{fullname}(\underline{anna} \underline{gialli}), \underline{status}), E \rangle, \langle \underline{cpage}(\underline{title}(\underline{Algebra})), E \rangle \}$$

Clearly, the fixpoint of  $\mathcal{J}_M$  (and hence  $Req_{M,W}$ ) for an arbitrary Web specification might be infinite. Consider for instance the following example.

*Example 13* Let  $W \equiv \{h(g(0), f(0))\}$  be a Web site and

$$I_M \equiv \{h(g(X)) \rightarrow h(g(g(X)))\langle q \rangle\}$$

be a set of completeness rules of a Web specification  $S$ . Then,

$$Req_{M,W} = \{ \langle h(g(g(0))), q \rangle, \langle h(g(g(g(0)))) \rangle, q \rangle, \langle h(g(g(g(g(0)))) \rangle, q \rangle, \dots \}$$

is an infinite set of requirements.

Fortunately, the computation of the set of requirements is finite for some interesting classes of Web specifications. Trivially, non-recursive specifications allow to reach  $lfp(\mathcal{J}_M)$  after a finite number of applications of  $\mathcal{J}_M$ , i.e.,  $lfp(\mathcal{J}_M) = \mathcal{J}_M \uparrow^W k, k \in \mathbb{N}$ . However, non-recursive definitions are not expressive enough for verification purposes, since some relevant conditions about Web sites cannot be formalized without resorting to recursion; e.g., some properties stated in Example 6 cannot be formulated by using a non-recursive specification.

In the following, we ascertain two important classes of recursive Web specifications whose set of requirements is finite. Basically, the idea is to consider those specifications in which the computation of the least fixpoint only generates marked expressions whose size is bounded [2]. We believe that such classes are expressive enough, since they can model the most common completeness as well as correctness requirements that one

can desire for a given Web site. Moreover, these classes allow us to generate only finite sets of conditions to be verified, which is *crucial* for the effectiveness of the whole method and *reasonable* in this context as Web sites are finite collections of semistructured data with an associated finite semantics.

The following definition formalizes the class of the *bounded* Web specifications.

**Definition 11** A Web specification  $(I_N \uplus I_M, R)$  is bounded iff, for each  $l \equiv (V_1, E_1, r_1, label_1) \in Lhs_M, r \equiv (V_2, E_2, r_2, label_2) \in RhS_M$  and each minimal injective simulation  $S$  of  $l$  in  $r|_V$  w.r.t.  $\sim, v \in V_2$ , the following properties hold

1. if  $v_2 \in \pi(S)$  and  $label_2(v_2) \in Var(r|_V)$ , then for all  $v_1 \in V_1$  s.t.  $label_1(v_1) \in Var(l)$ ,  $depth(r|_V, v_2) = depth(l, v_1)$ .
2. For each  $r \in RhS_M$ ,  $r$  does not contain any symbol of  $\Sigma_R$ .

Roughly speaking, Definition 11 states that, whenever the left-hand side  $l$  of a rule is simulated by (a subterm of) the right-hand side  $r$  of a (possibly different) rule, then no variables in the recognized substructure of  $r$  must be located at positions which are deeper than all the positions of the variables in  $l$ . Moreover, bounded Web specifications do not allow function calls in the completeness rules. Let us see an example.

*Example 14* Consider again the set of completeness rules  $I_M$  of Web specification  $S$  in Example 13. The left-hand side of the rule  $h(g(X)) \rightarrow h(g(g(X)))\langle q \rangle$  is simulated by its own right-hand side. Moreover, variable  $X$  in the right-hand side is located at depth 3, while the unique variable in the left-hand side is at depth 2. Thus,  $S$  is not bounded.

Now, take into account the Web specification  $S' \equiv (I'_N \uplus I'_M, \emptyset)$ , whose set of completeness rules is

$$I'_M \equiv \{ \langle m(n(X)) \rightarrow h(n(X), s(s(X)))\langle q' \rangle, \langle h(n(X)) \rightarrow m(n(X), t)\langle q'' \rangle \rangle, \quad q', q'' \in \{A, E\}. \}$$

Then,  $m(n(X))$  is simulated by  $m(n(X), t)$  and  $h(n(X))$  is simulated by the term  $h(n(X), s(s(X)))$ . In both cases, variables occurring in the substructures of the right-hand sides which are recognized by simulation and variables of the respective left-hand sides are located at the same depth. Therefore, the Web specification  $S'$  is bounded.

For bounded Web specifications, the least fixpoint of the operator  $\mathcal{J}_M$  is finite as stated by the next proposition. This provides an effective method for computing the set of requirements  $Req_{M,W}$  for this class of specifications.

**Proposition 3** Let  $(I_N \uplus I_M, R)$  be a bounded Web specification and  $\mathbb{W}$  be a Web site. Then, there exists  $k \in \mathbb{N}$  such that  $\text{lfp}(\mathcal{J}_M) = \mathcal{J}_M \uparrow^{\mathbb{W}} k$ .

Let us now introduce a more general class of Web specifications in which defined functions can be invoked in the rhs's of the completeness rules. In order to keep the size of the derived terms bounded, the key idea is to ensure that function calls do not generate infinite data structures.

Given a Web specification  $(I_N \uplus I_M, R)$ , we say that the completeness rule  $(1 \rightarrow \mu(x) \langle q \rangle) \in I_M$  is *function-dependent* w.r.t.  $R$  iff a function call  $f(t_1, \dots, t_n) \in \tau(\Sigma_R, \mathcal{V})$  occurs in  $x$ . Otherwise we say that rule  $r$  is *function-independent* w.r.t.  $R$ . Given a function-dependent rule  $r$ , we obtain the *function-independent version* of  $r$  w.r.t.  $R$ , and we denote it by  $r^*$ , by replacing all the function calls in the right-hand side of  $r$  with fresh variables not occurring in the rule.

*Example 15* Consider the function-dependent rule of Example 6

$$\text{member}(\text{name}(X), \text{surname}(Y)) \rightarrow \text{hpage}(\text{fullname}(\text{append}(X, Y), \text{status}) \langle E \rangle).$$

The function-independent version of the considered rule is

$$\text{member}(\text{name}(X), \text{surname}(Y)) \rightarrow \text{hpage}(\text{fullname}(Z), \text{status}) \langle E \rangle.$$

**Definition 12** A Web specification  $(I_N \uplus I_M, R)$  is bounded\* iff

1.  $(I_N \uplus \{r \in I_M \mid r \text{ is function-independent w.r.t. } R\}, R)$  is bounded.
2. For each function-dependent rule  $r \in I_M$  w.r.t.  $R$  and  $1 \in \text{Lhs}_M$ 
  - (a) There exists no minimal, injective simulation of  $1$  in any subterm of the right-hand side of  $r^*$  w.r.t.  $\sim$ .
  - (b) There exists no minimal, injective simulation of any subterm of the right-hand side of  $r^*$  in  $1$  w.r.t.  $\sim$ .
3.  $\{f \mid f(t_1, \dots, t_n) \rightarrow \mu(x) \in I_M\} \cap \Sigma_R = \emptyset$ .

Roughly speaking, by Definition 12, the rhs's of function-dependent completeness rules do not introduce terms which can be partially rewritten in a subsequent step. This suffices to ensure the finiteness of all partial rewriting sequences. Now, it is immediate to prove the finiteness of the least fixpoint of the  $\mathcal{J}_M$  operator for bounded\* Web specifications.

The following proposition generalizes Proposition 3 for bounded\* Web specifications.

**Proposition 4** Let  $(I_N \uplus I_M, R)$  be a bounded\* Web specification and  $\mathbb{W}$  be a Web site. Then, there exists  $k \in \mathbb{N}$  such that  $\text{lfp}(\mathcal{J}_M) = \mathcal{J}_M \uparrow^{\mathbb{W}} k$ .

Let us see an example.

*Example 16* Let us consider the Web specification of Example 6. Then, we can easily check that it is bounded\*. Moreover, the least fixpoint of  $\mathcal{J}_M$  is finite as witnessed by Example 12.

*Example 17* Consider now the following Web specification.

$$\begin{aligned} I_N &= \emptyset \\ I_M &= \{h(c(X)) \rightarrow h(g(X))\} \\ R &= \{g(0) \rightarrow c(c(0)), g(c(X)) \rightarrow c(g(X))\} \end{aligned}$$

The function-independent version of  $h(c(X)) \rightarrow h(g(X))$  is  $h(c(X)) \rightarrow h(Z)$ , and there exists a minimal, injective simulation of  $h(Z)$  in  $h(c(X))$  w.r.t.  $\sim$ . Therefore, the considered Web specification is not bounded\*. Actually, we can generate the following infinite partial rewrite sequence:

$$h(c(0)) \rightarrow h(c(c(0))) \rightarrow h(c(c(c(0)))) \dots$$

Now, we are ready to formalize step 2, that is, checking the computed completeness requirements in a given Web site. To accomplish this task, we first use simulation for checking whether (the marked part of) a requirement is embedded into some Web page of the considered site, and then consider the quantification attributes in order to diagnose completeness errors.

**Definition 13** Let  $\mathbb{W}$  be a Web site,  $(I_N \uplus I_M, R)$  be a Web specification and  $\text{Req}_{M, \mathbb{W}}$  be the set of requirements for  $\mathbb{W}$  w.r.t.  $I_M$ . Let  $\langle \mu(e), q \rangle \in \text{Req}_{M, \mathbb{W}}$ . The test set w.r.t.  $\langle \mu(e), q \rangle$  is defined as

$$\text{TEST}_{\langle \mu(e), q \rangle} = \{p \equiv (V, E, r, \text{label}) \in \mathbb{W} \mid \exists \text{ a minimal injective simulation of } \text{mark}(e, \mu) \text{ in } p|_V \text{ w.r.t. } \sim, \text{ with } v \in V\}.$$

Roughly speaking, this definition allows us to compute a subset of the Web site containing all the Web pages which simulate the marked part of a given requirement. These Web pages might be incomplete w.r.t. the Web specification, since they may not contain the considered requirement. Let us see an example.

*Example 18* Let us consider the completeness rule  $r$

$$\text{hpage}(\text{status}(\text{professor})) \rightarrow \text{hpage}(\text{status}(\text{professor}), \text{teaching})$$

and the Web site  $\mathbb{W}$  of Example 4. Rule  $r$  allows us to check whether Web pages of the professors contain

some teaching information. In order to do this, we use the marking information in the rhs of  $r$  to select the professor Web pages. Let us consider the requirement  $\langle \mu_1(e_1), A \rangle \equiv \langle \text{hpage}(\text{status}(\text{professor}, \text{teaching}), A) \rangle$ , which can be derived from  $W$  by means of  $r$ . By applying Definition 13, we get the following test set  $\text{TEST}_{\langle \mu_1(e_1), A \rangle}$

- { (2)  $\text{hpage}(\text{fullname}(\text{mariorossi}), \text{phone}(3333), \text{status}(\text{professor}), \text{hobbies}(\text{hobby}(\text{reading}), \text{hobby}(\text{gardening})))$ ,
- (4)  $\text{hpage}(\text{fullname}(\text{annagialli}), \text{status}(\text{professor}), \text{blink}(\text{phone}(4444)), \text{teaching}(\text{courselink}(\text{url}(\text{http://www.algebra.org}), \text{urlname}(\text{Algebra}))))$  }

which contains the two professor Web pages where the computed completeness requirement must be checked.

In the following, we consider completeness errors which refer to incomplete and/or missing Web pages. We distinguish two cases: the former allows us to discover whether a universal requirement is not fulfilled by a given Web site, while the latter recognizes unsatisfied existential requirements. In both cases, our analysis provides the missing/incomplete Web pages which are associated with those requirements.

**Definition 14** *Let  $W$  be a Web site,  $(I_N \uplus I_M, R)$  be a Web specification and  $\text{Req}_{M,W}$  be the set of requirements for  $W$  w.r.t.  $I_M$ . Let  $re \equiv \langle \mu(e), A \rangle \in \text{Req}_{M,W}$  be a universal requirement. Then,  $re$  is not satisfied in  $W$  if one of the following conditions hold:*

1.  $\text{TEST}_{\langle \mu(e), A \rangle} = \emptyset$ .
2. *There exists  $p \equiv (V, E, r, \text{label}) \in \text{TEST}_{\langle \mu(e), A \rangle}$  s.t. no minimal, injective simulation of  $e$  in  $p|_V$  w.r.t.  $\sim$ , with  $v \in V$ , exists.*

Vice versa, a universal requirement  $re$  is satisfied whenever it is possible to recognize  $re$  inside any Web page of the corresponding test set  $\text{TEST}_{re}$ . Let us clarify Definition 14 by an example.

**Example 19** Consider the Web site  $W$  of Example 4 and the universal requirement

$$\langle \mu_1(e_1), A \rangle \equiv \langle (\text{hpage}(\text{status}(\text{professor}, \text{teaching})), A) \rangle$$

belonging to the set of requirements  $\text{Req}_{M,W}$  of Example 12. The requirement simply states that *all* professor’s home pages must contain teaching information.

The test set associated with  $\langle \mu_1(e_1), A \rangle$ ,  $\text{TEST}_{\langle \mu_1(e_1), A \rangle}$ , was computed in Example 18 and contains all the

professor’s Web pages of the considered site. Now, by applying Definition 14, we detect that  $\langle \mu_1(e_1), A \rangle$  is not satisfied by  $W$ , since there does not exist any minimal, injective simulation of  $e_1$  in (a subterm of) Web page (2) w.r.t.  $\sim$ .

In fact, Web page (2) lacks teaching information.

Finally, an existential requirement  $re$  is fulfilled, if it is recognized inside (at least) a Web page which belongs to the test set  $\text{TEST}_{re}$ .

**Definition 15** *Let  $W$  be a Web site,  $(I_N \uplus I_M, R)$  be a Web specification and  $\text{Req}_{M,W}$  be the set of requirements for  $W$  w.r.t.  $I_M$ . Let  $re \equiv \langle \mu(e), E \rangle \in \text{Req}_{M,W}$  be an existential requirement. Then,  $re$  is not satisfied in  $W$  if one of the following conditions hold:*

1.  $\text{TEST}_{\langle \mu(e), E \rangle} = \emptyset$ .
2. *For each  $p \equiv (V, E, r, \text{label}) \in \text{TEST}_{\langle \mu(e), E \rangle}$  no minimal, injective simulation of  $e$  in  $p|_V$  w.r.t.  $\sim$ , with  $v \in V$ , exists.*

**Example 20** Consider the Web site  $W$  of Example 4 and the existential requirement

$$\langle \mu_2(e_2), E \rangle \equiv \langle \text{cpage}(\text{title}(\text{Algebra})), E \rangle$$

belonging to the set of requirements  $\text{Req}_{M,W}$  of Example 12. Since  $\text{TEST}_{\langle \mu_2(e_2), E \rangle}$  is empty, by Definition 15,  $re$  is not satisfied in  $W$ . This implies that a Web page containing some information about the Algebra course should be provided.

Consider now the following existential requirements in  $\text{Req}_{M,W}$

$$\langle \mu_3(e_3), E \rangle \equiv \langle \text{member}(\text{name}(\text{anna}), \text{surname}(\text{gialli})), E \rangle$$

$$\langle \mu_4(e_4), E \rangle \equiv \langle \text{hpage}(\text{fullname}(\text{giulioverdi}), \text{status}), E \rangle$$

We have that  $\text{TEST}_{\langle \mu_3(e_3), E \rangle}$  just includes Web page (1) of  $W$ , while  $\text{TEST}_{\langle \mu_4(e_4), E \rangle}$  contains Web pages (2), (3) and (4). For both requirements, there is no Web page  $p$  in the test sets such that a minimal, injective simulation of  $e_3$  (respectively,  $e_4$ ) in (a subterm of)  $p$  w.r.t.  $\sim$  exists. Therefore, the Web site  $W$  does not meet the given requirements. More precisely, the former unsatisfied requirement states that an entry for *Anna Gialli* must be introduced in the group member list, and the latter detects that the home page of *Giulio Verdi* is missing.

The requirements which are not fulfilled can be considered as *incompleteness symptoms*. This allows us not only to locate bugs and inconsistencies w.r.t. a given specification, but also to easily repair them by comparing incomplete pages to unsatisfied requirements, since the latter ones provide the missing information which is needed to complete the erroneous Web pages.

An algorithm for detecting incomplete information can be defined as follows. First of all, we generate the requirements to be checked by computing the fixpoint of the operator  $\mathcal{J}_M$ . Next, for each (universal/existential) requirement, the corresponding test set is produced and then checked in order to analyze which requirements are not fulfilled. The analysis is based on Definition 14 (resp., Definition 15) which formalizes the notion of universal (resp., existential) requirement satisfiability. Finally, incompleteness symptoms are returned, which are needed to help the user to locate the missing/incomplete pages, whenever a requirement is not verified.

The basic procedure is sketched in Algorithm 2, which takes as input a Web site  $W$ , a set of completeness rules  $I_M$ , and a canonical TRS  $R$ .

**Algorithm 2** An algorithm for detecting completeness errors.

```

1: procedure COMPLETENESS-ERRORS( $W, I_M, R$ )
2:    $Req_{M,W} \leftarrow lfp(\mathcal{J}_M) \setminus W$ 
3:   for all  $\langle \mu(e), q \rangle \in Req_{M,W}$  do
4:      $TEST_{\langle \mu(e), q \rangle} \leftarrow \{p \equiv (V, E, r, label) \in W \mid$ 
        $mark(e, \mu) \cong p|_V, v \in V\}$ 
5:     if  $TEST_{\langle \mu(e), q \rangle} = \emptyset$  then
6:       output("Error: A Web page containing  $e$  must occur
         in Web site  $W$ !")
7:     end if
8:     case  $q$ 
9:        $q = A$  :
10:      if  $\exists p \equiv (V, E, r, label) \in TEST_{\langle \mu(e), A \rangle}$  s.t.  $e \not\cong p|_V$ ,
        with  $v \in V$  then
11:        output("Error:  $e$  must occur in all Web pages
          of  $TEST_{\langle \mu(e), A \rangle}$ !")
12:      end if
13:       $q = E$  :
14:      if  $\forall p \equiv (V, E, r, label) \in TEST_{\langle \mu(e), E \rangle}$ ,  $e \not\cong p|_V$ ,
        with  $v \in V$  then
15:        output("Error:  $e$  must occur in at least one
          Web page of  $TEST_{\langle \mu(e), E \rangle}$ !")
16:      end if
17:    end case
18:  end for
19: end procedure

```

**Proposition 5** Let  $W$  be a Web site,  $(I_N \uplus I_M, R)$  be a bounded\* Web specification, and  $Req_{M,W}$  be the set of requirements for  $W$  w.r.t.  $I_M$ . Then, the procedure COMPLETENESS-ERRORS( $W, I_M, R$ ) terminates. Moreover, for each error message regarding term  $e$ , which is returned by the procedure, there exists  $\langle \mu(e), q \rangle \in Req_{M,W}$ , with  $q \in \{A, E\}$ , which is not satisfied in  $W$ .

**7 The GVERDI system**

The basic methodology presented so far has been implemented in the prototype system GVERDI (VERification

and Rewriting for Debugging Internet sites) which has been written in Haskell (GHC v6.2.2) and is publicly available together with a set of tests at

<http://www.dsic.upv.es/users/elp/GVerdi/>

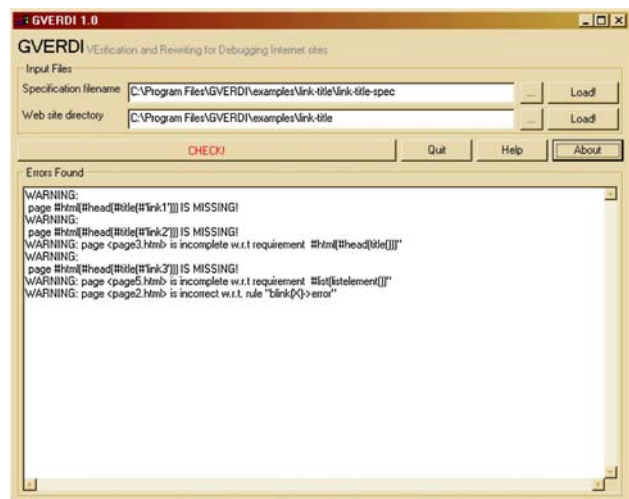
A short system description can be found in [3,6], while a thorough user’s manual of the tool can be retrieved at the URL mentioned above.

The implementation consists of approximately 1,100 lines of source code. It includes a parser for semistructured expressions (i.e., XML/XHTML documents) and Web specifications, and several modules implementing the partial rewriting mechanism, the verification technique, and the graphical user interface. A snapshot of the running system is shown in Fig. 8.

The system allows the user to load a Web site together with a Web specification. Additionally, he/she can inspect the loaded data and finally check the Web pages w.r.t. the Web site specification. We have tested the system on several XHTML Web sites and XML data collections which can be found at the URL address mentioned above. For instance, we checked the Web site of the Computational Logic Group of the University of Udine which is available at <http://www.dimi.uni-ud.it/clg>. It contains about 20 Web pages concerning publications, people, and projects of the group. In each considered test case, we have been able to detect the errors (i.e., missing and incomplete Web pages) efficiently.

**8 Conclusions**

Conceiving and maintaining Web sites is a difficult task. In this paper, we provide a rewriting-based, formal



**Fig. 8** A screenshot of the system

specification language which can be used to impose properties both on the structure (syntactic properties) and on the contents (semantic properties) of Web sites. Some XML schema languages such as XLINKIT [19] and Schematron [35] can express some of the semantic constraints we consider; however, our rule specification language, which offers the expressiveness and computational power of functions, is richer and can be appreciated much by users who prefer to avoid the encumbrances of DTDs and XML rule languages. The computation mechanism underlying our language is based on a novel rewriting-like technique, called *partial rewriting*, in which the traditional pattern matching mechanism is replaced by tree *simulation* [24]. In our methodology, Web sites are automatically checked w.r.t. a given Web specification in order to detect incorrect, incomplete, and missing Web pages. Moreover, by analyzing the error symptoms, we are also able to find out the missing information needed to repair the Web site and the wrong information we have to remove or to fix. We have also discussed some implementation details of the system GVERDI, a prototype implementation of the verification framework which can help Web administrators to design, check, and maintain their Web sites.

Let us conclude by mentioning some directions for future work. We are currently extending our framework in order to provide a method for synthesizing the marking information semi-automatically (currently, marks are provided by the user). Besides, we plan to extend the GVERDI system with a graphical, high-level language for the specification of the conditions to be checked in order to make the system easy to use and intuitive even for the users who have no expertise in formal methods. Finally, in order to consider semistructured documents containing cycles, which could be easily generated —for instance— by using ID and ID-REF XML attributes, a graph rewriting extension of our framework is also envisaged. In fact, trees are not enough to model loops, we believe that graph structures and the associated rewriting framework are essential to deal with such circular documents and to enhance the expressiveness of our specification language.

**Acknowledgments** We gratefully thank Javier García-Vivó for his help with the implementation of GVERDI and the experiments with the system. This work has been partially supported by the EU (FEDER) and the Spanish MEC, under grant TIN 2004-7943-C04-02, by ICT for EU-India Cross Cultural Dissemination Project under grant ALA/95/23/2003/077-054, and by Generalitat Valenciana under grant GR03/025.

## Appendix: Proofs of the technical results

In this appendix we provide the proofs of the main technical results of the paper.

**Proposition 1** Let  $s_1 \equiv (V_1, E_1, r_1, \text{label}_1)$ ,  $s_2 \equiv (V_2, E_2, r_2, \text{label}_2)$  be two Web page templates in  $\tau(\text{Text} \cup \text{Tag}, \mathcal{V})$ . A minimal page simulation  $\mathbf{S}$  of  $s_1$  in  $s_2$  w.r.t.  $\sim$  is a mapping  $\mathbf{S} : V_1 \rightarrow V_2$ .

*Proof* Let  $s_1 \equiv (V_1, E_1, r_1, \text{label}_1)$ ,  $s_2 \equiv (V_2, E_2, r_2, \text{label}_2) \in \tau(\text{Text} \cup \text{Tag}, \mathcal{V})$ . By Conditions 1 and 3 of Definition 5, and by using the fact that  $s_1$  has an underlying tree structure (in particular, it is a connected graph), we have that

$\forall v_1 \in V_1, \exists v_2 \in V_2$  such that  $(v_1, v_2) \in \mathbf{S}$ .

Moreover, the minimality of  $\mathbf{S}$  ensures that

$\forall v_1 \in V_1, \exists$  a *unique*  $v_2 \in V_2$  such that  $(v_1, v_2) \in \mathbf{S}$

which implies that  $\mathbf{S}$  is a mapping from  $V_1$  to  $V_2$ .

**Proposition 2** Let  $\mathbb{W}$  be a Web site, and  $(\mathcal{I}_N \uplus \mathcal{I}_M, R)$  be a Web specification. Then, the procedure CORRECTNESS-ERRORS( $\mathbb{W}, \mathcal{I}_N, R$ ) terminates, and for each pair  $(p, \text{ls})$ , which is returned,  $p$  is an incorrect Web page w.r.t.  $(\mathcal{I}_N \uplus \mathcal{I}_M, R)$  and  $\text{ls}$  is the corresponding incorrectness symptom.

*Proof* First of all, let us prove the termination of the algorithm. The outer loop (lines 2–9) is executed  $|\mathbb{W}|$  times. The inner loop (lines 3–8) is executed  $|\mathcal{I}_N|$  times. As for the partial rewrite step  $p \xrightarrow{\sigma}_{r_u}$  error in line 4, it terminates, because it is just an application of the simulation algorithm [24], which terminates. The evaluation of the condition in line 5 terminates, since:

1. The problem of evaluating  $X_i \sigma \in \mathcal{L}(\text{rexp}_i)$ ,  $i = 1, \dots, n$  boils down to the membership problem in regular languages, which is decidable (see [25]), so there exists an effective procedure which tests  $(X_i \sigma \in \mathcal{L}(\text{rexp}_i), i = 1, \dots, n)$  in finite time.
2. Evaluating the condition  $[(s_j = t_j)\sigma, j = 1, \dots, m, \text{ holds in } R]$  trivially terminates, since  $R$  is canonical (in particular, it is terminating).

The output command in line 6 clearly terminates. Summing up, we execute the block of terminating instructions, which is included in lines 4–7,  $|\mathbb{W}| * |\mathcal{I}_N|$  times, so the whole procedure terminates.

Let us prove the partial correctness of the algorithm, that is, if the procedure terminates producing the outcome  $(p, \text{ls})$ , then  $p \in \mathbb{W}$  is incorrect w.r.t.  $(\mathcal{I}_N \uplus \mathcal{I}_M, R)$  and  $\text{ls}$  is an incorrect symptom for  $p$ .

Let  $(p, \text{ls})$  be an output of the procedure. This implies that (1)  $p \xrightarrow{\sigma}_{r_u}$  error for some  $r \equiv (1 \rightarrow \text{error} \mid X_1 \text{ in rexp}_1, \dots, X_n \text{ in rexp}_n, s_1 = t_1, \dots, s_m = t_m) \in \mathcal{I}_N$  (see line 4); (2) for each  $i = 1, \dots, n$ ,  $(X_i \sigma \in \mathcal{L}(\text{rexp}_i), i = 1, \dots, n)$  (see line 5); and (3) for each  $j = 1, \dots, m$ ,  $(s_j = t_j)\sigma$  holds in  $R$  (see line 5). Now, by simply applying Definition 9, we get that  $p \in \mathbb{W}$  is incorrect w.r.t.  $(\mathcal{I}_N \uplus \mathcal{I}_M, R)$  and  $\text{ls}$  is an incorrect symptom for  $p$ .



In order to prove Proposition 3, we need the following auxiliary definitions and results.

**Definition 16** Given a term  $t \in \tau(\Sigma, \mathcal{V})$ , the height of  $t$ ,  $\text{height}(t)$ , is defined as follows.

$$\text{height}(t) = \begin{cases} 0 & \text{if } t \equiv X \in \mathcal{V} \text{ or } t \equiv c \in \tau(\Sigma, \mathcal{V}) \\ 1 + \max\{\text{height}(t_i) \mid i = 1, \dots, n\} & \text{if } t \equiv f(t_1, \dots, t_n) \in \tau(\Sigma, \mathcal{V}) \end{cases}$$

We can lift the notion of height to substitutions in the following way.

**Definition 17** Given a substitution

$\sigma = \{X_1/t_1, \dots, X_n/t_n\}$ , the height of  $\sigma$ ,  $\text{height}(\sigma)$ , is defined as follows.

$$\text{height}(\sigma) = \max\{\text{height}(t_i) \mid i = 1, \dots, n\}$$

Now we can prove two useful technical results.

**Proposition 6** Let  $(\mathbb{I}_N \uplus \mathbb{I}_M, R)$  be a bounded Web specification and

$$\mu_0(s_0) \xrightarrow{r_0^{\sigma_0}} \mu_1(s_1) \xrightarrow{r_1^{\sigma_1}} \mu_2(s_2) \xrightarrow{r_2^{\sigma_2}} \dots$$

be a partial rewrite sequence, where  $r_i \in \mathbb{I}_M$ ,  $i = 0, 1, 2, \dots$ . Then, for each  $s_i$ ,  $i = 1, 2, \dots$ ,

$$s_i \equiv r_{i-1}\sigma_{i-1}$$

where  $r_{i-1}$  is the right-hand side of the rule  $r_{i-1}$ .

*Proof* It directly comes from Definition 8 and by the fact that  $\text{Reduce}(r\sigma, R) = r\sigma$  when considering bounded Web specification, since no function of  $R$  can appear in the right-hand sides of the completeness rules.

In other words, Proposition 6 states that each term occurring in a partial rewrite sequence is an instance of the right-hand side of some completeness rule.

**Proposition 7** Let  $(\mathbb{I}_N \uplus \mathbb{I}_M, R)$  be a bounded Web specification and

$$\mu_0(s_0) \xrightarrow{r_0^{\sigma_0}} \mu_1(s_1) \xrightarrow{r_1^{\sigma_1}} \mu_2(s_2) \xrightarrow{r_2^{\sigma_2}} \dots$$

be a partial rewrite sequence, where  $r_i \in \mathbb{I}_M$ ,  $i = 0, 1, 2, \dots$ . Then, for each  $s_i$ ,  $i = 1, 2, \dots$ ,

$$\text{height}(s_i) \leq \text{height}(r_{i-1}) + \text{height}(\sigma_{i-1})$$

where  $r_{i-1}$  is the right-hand side of the rule  $r_{i-1}$ .

*Proof* It is a direct consequence of Proposition 6 Definition 16 and Definition 17.

**Lemma 1** Let  $(\mathbb{I}_N \uplus \mathbb{I}_M, R)$  be a bounded Web specification and

$$\mu_0(s_0) \xrightarrow{r_0^{\sigma_0}} \mu_1(s_1) \xrightarrow{r_1^{\sigma_1}} \mu_2(s_2) \xrightarrow{r_2^{\sigma_2}} \dots$$

be a (possibly infinite) partial rewrite sequence, where  $r_i \in \mathbb{I}_M$ ,  $i = 0, 1, 2, \dots$ . Then, for each  $s_i$ ,  $i = 1, 2, \dots$ ,

$$\text{height}(s_i) \leq \text{height}(r_{i-1}) + \text{height}(\sigma_0)$$

where  $r_{i-1}$  is the right-hand side of rule  $r_{i-1}$ .

*Proof* By contradiction, let  $s_{k'}$  be the first term appearing in the partial rewrite sequence such that

$$\text{height}(s_{k'}) > \text{height}(r_{k'-1}) + \text{height}(\sigma_0),$$

where  $r_{k'-1}$  is the right-hand side of the rule  $r_{k'-1}$ . Clearly,  $k' > 1$ , as the claim trivially holds for  $k' = 1$ . Moreover, for each  $j = 1, \dots, k' - 1$ , we have

$$\text{height}(s_j) \leq \text{height}(r_{j-1}) + \text{height}(\sigma_0).$$

Now, let us focus on the partial rewrite step

$$s_{k'-1} \xrightarrow{r_{k'-1}^{\sigma_{k'-1}}} s_{k'}.$$

By Proposition 6,  $s_{k'-1} \equiv r_{k'-2}\sigma_{k'-2}$ . So, we have

$$\text{height}(s_{k'-1}) \leq \text{height}(r_{k'-2}) + \text{height}(\sigma_{k'-2}).$$

And, also  $\text{height}(\sigma_{k'-2}) \leq \text{height}(\sigma_0)$ .

Again, by Proposition 6,  $s_{k'} \equiv r_{k'-1}\sigma_{k'-1}$ . And, hence,

$$\text{height}(s_{k'}) \leq \text{height}(r_{k'-1}) + \text{height}(\sigma_{k'-1}) \quad \text{by Proposition 7.}$$

Here, we distinguish two cases.

Since  $s_{k'-1} \equiv r_{k'-2}\sigma_{k'-2}$ , the partial rewrite step  $s_{k'-1} \xrightarrow{r_{k'-1}^{\sigma_{k'-1}}} s_{k'}$  can be given on a vertex in  $r_{k'-1}$  (Case 1) or on a vertex belonging to a term in  $\{t \mid X/t \in \sigma_{k'-2}\}$  (Case 2).

Case 1. Since the Web specification is bounded, no variables in the substructure of  $r_{k'-1}$ , which is simulated by the left-hand side of rule  $r_{k'-1}$ , can be located at positions which are deeper than all the positions of the variables in the considered left-hand side. So, we must have  $\text{height}(\sigma_{k'-1}) = \text{height}(\sigma_{k'-2})$ . Therefore,

$$\begin{aligned} \text{height}(s_{k'}) &\leq \text{height}(r_{k'-1}) + \text{height}(\sigma_{k'-1}) \\ &= \text{height}(r_{k'-1}) + \text{height}(\sigma_{k'-2}) \\ &\leq \text{height}(r_{k'-1}) + \text{height}(\sigma_0) \end{aligned}$$

which is a contradiction.

Case 2. Trivially,  $\text{height}(\sigma_{k'-1}) \leq \text{height}(\sigma_{k'-2})$ .

Hence,

$$\begin{aligned} \text{height}(s_{k'}) &\leq \text{height}(r_{k'-1}) + \text{height}(\sigma_{k'-1}) \\ &\leq \text{height}(r_{k'-1}) + \text{height}(\sigma_{k'-2}) \\ &\leq \text{height}(r_{k'-1}) + \text{height}(\sigma_0) \end{aligned}$$

which also leads to a contradiction.

Thus, there cannot exist a term  $s_{k'}$  such that  $\text{height}(s_{k'}) > \text{height}(r_{k'-1}) + \text{height}(\sigma_0)$  and the claim is proven.

**Proposition 8** Let  $\mathbb{W}$  be a Web site,  $(\mathbb{I}_N \uplus \mathbb{I}_M, R)$  be a Web specification. Then,

$$\begin{aligned} \{\mu(s) \mid \varepsilon(p) \xrightarrow{m} \mu(s), 0 \leq m \leq k, \varepsilon(p) \in \mathbb{W}\} \\ = \{\mu(s) \mid \langle \mu(s), \alpha \rangle \in \mathcal{J}_M \uparrow^{\mathbb{W}} k\} \end{aligned}$$

*Proof*

- ( $\supseteq$ ) Simple induction on the ordinal power  $k$ .
- ( $\subseteq$ ) Simple induction on the length  $m$  of the rewrite sequence.

**Corollary 1** *Let  $\mathbb{W}$  be a Web site,  $(\mathbb{I}_N \uplus \mathbb{I}_M, R)$  be a Web specification. Then,*

$$\{\mu(s) \mid \varepsilon(p) \xrightarrow{*} \mu(s), \varepsilon(p) \in \mathbb{W}\} = \{\mu(s) \mid \langle \mu(s), \mathfrak{q} \rangle \in \text{lfp}(\mathcal{J}_M)\}$$

Now, we are ready to prove Proposition 3.

**Proposition 3** *Let  $(\mathbb{I}_N \uplus \mathbb{I}_M, R)$  be a bounded Web specification and  $\mathbb{W}$  be a Web site. Then, there exists  $k \in \mathbb{N}$  such that  $\text{lfp}(\mathcal{J}_M) = \mathcal{J}_M \uparrow^{\mathbb{W}} k$ .*

*Proof* First, let us show that each  $\mu(s)$ , s.t.  $\langle \mu(s), \mathfrak{q} \rangle \in \text{lfp}(\mathcal{J}_M)$ , has a bounded height.

By Corollary 1,  $\{\mu(s) \mid \varepsilon(p) \xrightarrow{*} \mu(s), \varepsilon(p) \in \mathbb{W}\} = \{\mu(s) \mid \langle \mu(s), \mathfrak{q} \rangle \in \text{lfp}(\mathcal{J}_M)\}$ .

Thus, for each  $\mu(s)$ , such that  $\langle \mu(s), \mathfrak{q} \rangle \in \text{lfp}(\mathcal{J}_M)$ , there exists a Web page  $\varepsilon(p) \in \mathbb{W}$  such that  $\varepsilon(p) \xrightarrow{*} \mu(s)$ . Here, we distinguish two cases.

Case  $(\varepsilon(p) \xrightarrow{0} \mu(s))$ . Let

$$H_W = \max\{\text{height}(p) \mid \varepsilon(p) \in \mathbb{W}\}.$$

In this case,  $\mu(s) \equiv \varepsilon(p)$ , hence

$$\text{height}(\mu(s)) = \text{height}(\varepsilon(p)) \leq H_W.$$

Case  $(\varepsilon(p) \xrightarrow{+} \mu(s))$ . Let

$$H_I = \max\{\text{height}(r) \mid r \in \text{Rhs}_M\}$$

$$H_S = \max\{\text{height}(\sigma') \mid \varepsilon(p) \xrightarrow{\sigma'} \mu(s), \varepsilon(p) \in \mathbb{W}, r' \in \mathbb{I}_M\}.$$

Now, we have  $\varepsilon(p) \xrightarrow{\sigma_0} \mu(s_1) \xrightarrow{*} \mu(s)$ . By Lemma 1, there exists  $r \in \text{Rhs}_M$  such that

$$\text{height}(s) \leq \text{height}(r) + \text{height}(\sigma_0) \leq H_I + H_S.$$

Thus, for each  $\mu(s)$ , such that  $\langle \mu(s), \mathfrak{q} \rangle \in \text{lfp}(\mathcal{J}_M)$ , we get

$$\text{height}(\mu(s)) \leq \max\{H_W, H_I + H_S\}.$$

Since we have only a finite number of marked terms whose height is less than or equal to  $\max\{H_W, H_I + H_S\}$ ,  $\text{lfp}(\mathcal{J}_M)$  must be a finite set. Finally, since the operator  $\mathcal{J}_M$  is continuous (in particular, monotonic) and  $\text{lfp}(\mathcal{J}_M)$  is a finite set, then there exists a natural number  $k$  such that  $\text{lfp}(\mathcal{J}_M) = \mathcal{J}_M \uparrow^{\mathbb{W}} k$ .

In order to prove Proposition 4, we first give some auxiliary results.

**Lemma 2** *Let  $s_1$  and  $s_2$  be two Web page templates in  $\tau(\text{Text} \cup \text{Tag}, \mathcal{V})$ . Then,  $s_1 \cong s_2$  and  $s_2 \cong s_1$  iff there exists a substitution  $\sigma$*

*$s_2\sigma$  partially matches  $s_1$  via some substitution  $\sigma'$ .*

*Proof*

( $\Rightarrow$ ) Let us assume that

$$s_1 \cong s_2 \text{ or } s_2 \cong s_1.$$

Case  $s_1 \cong s_2$ . By Definition 5, we also have  $s_1 \cong s_2\sigma$ , for every  $\sigma$ . Directly, by Definition 7, there exists a substitution  $\sigma'$  s.t.  $s_2\sigma$  partially matches  $s_1$  via  $\sigma'$ .

Case  $s_2 \cong s_1$ .  $s_1$  partially matches  $s_2$  via  $\sigma''$ . By Definition 5, we also have  $s_2\sigma'' \cong s_1$ . Indeed,  $s_2\sigma''$  is embedded into  $s_1$ . Therefore,  $s_1 \cong s_2\sigma''$ , and  $s_2\sigma''$  partially matches  $s_1$  via the empty substitution  $\epsilon$ . So, if we take  $\sigma' = \epsilon$  and  $\sigma = \sigma''$ , then there exists  $\sigma$  and  $\sigma'$  such that  $s_2\sigma$  partially matches  $s_1$  via  $\sigma'$ .

( $\Leftarrow$ ) Since there exists a substitution  $\sigma$  s.t.  $s_2\sigma$  partially matches  $s_1$  via some substitution  $\sigma'$ , we have that  $s_1 \cong s_2\sigma$ . Here we can distinguish two cases: (1) there is a minimal, injective simulation of  $s_1$  in a subterm of  $s_2$  w.r.t.  $\sim$  and hence  $s_1 \cong s_2$ ; (2) there is a minimal, injective simulation of  $s_1$  in  $s_2\sigma$  w.r.t.  $\sim$ , which also involves some terms occurring in  $\sigma$ . In this case, we have that, for some  $X/t \in \sigma$ ,  $t$  is simulated by some subterm of  $s_1$ . It is not difficult to see that, by Definitions 4 and 5, there always exists a minimal, injective simulation between a variable and an arbitrary term. Exploiting this fact, we can state that each variable in  $s_2$ , which is replaced by some term in  $s_2\sigma$ , simulates the corresponding subterm in  $s_1$ , thus  $s_2 \cong s_1$ .

**Lemma 3** *Let  $(\mathbb{I}_N \uplus \mathbb{I}_M, R)$  be a bounded\* Web specification and  $r \equiv (1 \rightarrow r(\mathfrak{q})) \in \mathbb{I}_M$  be a function-dependent rule w.r.t.  $R$ . Let  $\mu_1(s_1)$  and  $\mu_2(s_2)$  be two marked Web page templates. If  $\mu_1(s_1) \xrightarrow{\sigma} \mu_2(s_2)$ , then there do not exist a Web page template  $\mu_3(s_3)$  and substitution  $\sigma'$  s.t.  $\mu_2(s_2) \xrightarrow{\sigma'} \mu_3(s_3)$ , with  $r' \in \mathbb{I}_M$ .*

*Proof* Let us consider  $\mu_1(s_1) \xrightarrow{\sigma} \mu_2(s_2)$ . Then, by Definition 8,

$$s_2 = \text{Reduce}(r\sigma, R).$$

We consider two cases.

Case  $s_2 \in \tau(\Sigma_R, \mathcal{V})$ .  $s_2 = \text{Reduce}(r\sigma, R)$  is just an irreducible form computed by the canonical TRS  $R$ . No partial rewriting steps can be applied to  $s_2$ , because (1)

$$\{\mathfrak{f} \mid \mathfrak{f}(t_1, \dots, t_n) \rightarrow \mu(r) \in \mathbb{I}_M\} \cap \Sigma_R = \emptyset,$$

since  $(\mathbb{I}_N \uplus \mathbb{I}_M, R)$  is bounded\*; (2) the function calls cannot be applied to arguments containing symbols in

$$\{\mathfrak{f} \mid \mathfrak{f}(t_1, \dots, t_n) \rightarrow \mu(r) \in \mathbb{I}_M\}.$$

Case  $s_2 \in \tau(\text{Text} \cup \text{Tag} \cup \Sigma_R, \mathcal{V}) \setminus \tau(\Sigma_R, \mathcal{V})$ .  
 $s_2 = \text{Reduce}(\tau\sigma, R) = \tau^*\sigma^*$ , where  $\tau^*$  is the right-hand side of the function-independent version of  $r$  w.r.t.  $R$ , and  $\sigma^*$  is a suitable substitution. Since  $(\mathbb{I}_N \uplus \mathbb{I}_M, R)$  is bounded\*, condition 2 of Definition 12 holds. Thus, by applying Lemma 2, we can ensure that the subterms that occur at positions coming from  $\tau^*$  cannot be partially rewritten in a subsequent step. Moreover, for each  $X/t \in \sigma^*$ ,  $t$  does not contain any subterm which can be partially rewritten, because (1)

$$\{f \mid f(\tau_1, \dots, \tau_n) \rightarrow \mu(x) \in \mathbb{I}_m\} \cap \Sigma_R = \emptyset,$$

since  $(\mathbb{I}_N \uplus \mathbb{I}_M, R)$  is bounded\*; (2) the function calls are never applied to arguments containing symbols in  $\{f \mid f(\tau_1, \dots, \tau_n) \rightarrow \mu(x) \in \mathbb{I}_m\}$ . Therefore, there are no subterms in  $s_2 = \tau^*\sigma^*$  which can be partially rewritten any longer.

**Lemma 4** *Let  $(\mathbb{I}_N \uplus \mathbb{I}_M, R)$  be a bounded\* Web specification and*

$$\mu_0(s_0) \xrightarrow{\sigma_0} \mu_1(s_1) \xrightarrow{\sigma_1} \mu_2(s_2) \xrightarrow{\sigma_2} \dots$$

*be a (possibly infinite) partial rewrite sequence, where  $r_i \in \mathbb{I}_m, i = 0, 1, 2, \dots$ . Then, there exists  $k \in \mathbb{N}$  s.t. for each  $s_i, i = 0, 1, 2, \dots$ ,*

$$\text{height}(s_i) \leq k.$$

*Proof* Let us consider the partial rewrite sequence

$$\mu_0(s_0) \xrightarrow{\sigma_0} \mu_1(s_1) \xrightarrow{\sigma_1} \mu_2(s_2) \xrightarrow{\sigma_2} \dots$$

We distinguish two cases.

Case 1. First consider that the partial rewrite sequence only contains function-independent rules. Since,  $(\mathbb{I}_N \uplus \{r \in \mathbb{I}_m \mid r \text{ is function-independent w.r.t. } R\}, R)$  is bounded, we can simply apply Lemma 1. So,

$$\text{height}(s_i) \leq \text{height}(\tau_{i-1}) + \text{height}(\sigma_0)$$

where  $\tau_{i-1}$  is the right-hand side of rule  $r_{i-1}, i = 1, 2, \dots$ . Let  $k' = \max\{\text{height}(\tau) \mid \tau \in \text{Rh}_{\Sigma_M}\}$ . By taking  $k = \max\{\text{height}(s_0), k' + \text{height}(\sigma_0)\}$ , we prove the claim.

Case 2. Now, assume that function-dependent completeness rules are also used in the sequence. Let us suppose that  $\mu_j(s_j) \xrightarrow{\sigma_j} \mu_{j+1}(s_{j+1})$  is the first partial rewrite step of the sequence in which a function-dependent rule  $r_i$  occurs. By Lemma 3,  $\mu_{j+1}(s_{j+1})$  cannot be partially rewritten any longer, therefore the sequence must be finite. So, by inspecting the heights of every term appearing in the sequence, we can find the one with the maximum height  $h$ . Finally, it is enough to choose  $k = h$  and the claim is proven.

**Proposition 4** *Let  $(\mathbb{I}_N \uplus \mathbb{I}_M, R)$  be a bounded\* Web specification and  $\mathbb{W}$  be a Web site. Then, there exists  $k \in \mathbb{N}$  such that  $\text{lfp}(\mathcal{J}_M) = \mathcal{J}_M \uparrow^{\mathbb{W}} k$ .*

*Proof* Similarly to the proof of Proposition 3, by Corollary 1,  $\{\mu(s) \mid \varepsilon(p) \xrightarrow{*} \mu(s), \varepsilon(p) \in \mathbb{W}\} = \{\mu(s) \mid \langle \mu(s), \mathfrak{q} \rangle \in \text{lfp}(\mathcal{J}_M)\}$ .

Thus, for each  $\mu(s)$ , such that  $\langle \mu(s), \mathfrak{q} \rangle \in \text{lfp}(\mathcal{J}_M)$ , there exists a partial rewrite sequence  $\varepsilon(p) \xrightarrow{*} \mu(s)$ , where  $\varepsilon(p) \in \mathbb{W}$ . As  $(\mathbb{I}_N \uplus \mathbb{I}_M, R)$  is bounded\*, every marked term appearing in the partial rewrite sequence  $\varepsilon(p) \xrightarrow{*} \mu(s)$  has a bounded height by Lemma 4 and hence  $\text{lfp}(\mathcal{J}_M)$  is a finite set.

Finally, since the operator  $\mathcal{J}_M$  is continuous (in particular, monotonic) and the set  $\text{lfp}(\mathcal{J}_M)$  is finite, then there exists a natural number  $k$  such that  $\text{lfp}(\mathcal{J}_M) = \mathcal{J}_M \uparrow^{\mathbb{W}} k$ .

**Proposition 5** *Let  $\mathbb{W}$  be a Web site,  $(\mathbb{I}_N \uplus \mathbb{I}_M, R)$  be a bounded\* Web specification, and  $\text{Req}_{\mathbb{Q}_M, \mathbb{W}}$  be the set of requirements for  $\mathbb{W}$  w.r.t.  $\mathbb{I}_M$ . Then, the procedure  $\text{COMPLETENESS-ERRORS}(\mathbb{W}, \mathbb{I}_M, R)$  terminates. Moreover, for each error message regarding term  $e$ , which is returned by the procedure, there exists  $\langle \mu(e), \mathfrak{q} \rangle \in \text{Req}_{\mathbb{Q}_M, \mathbb{W}}$ , with  $\mathfrak{q} \in \{A, E\}$ , which is not satisfied in  $\mathbb{W}$ .*

*Proof* First of all, let us prove the termination of procedure  $\text{COMPLETENESS-ERRORS}(\mathbb{W}, \mathbb{I}_M, R)$ . Line 2 computes the sets of requirements  $\text{Req}_{\mathbb{Q}_M, \mathbb{W}}$ . Since  $\mathbb{W}$  is a bounded\* Web specification, by Proposition 4, the  $\text{lfp}(\mathcal{J}_M)$  is a finite set which is computed in finite time. Therefore,  $\text{Req}_{\mathbb{Q}_M, \mathbb{W}} \leftarrow \text{lfp}(\mathcal{J}_M) \setminus \mathbb{W}$  is also computed in finite time.

The loop within lines 3–18 is executed  $|\text{Req}_{\mathbb{Q}_M, \mathbb{W}}|$  times. It simply takes every single requirement belonging to  $\text{Req}_{\mathbb{Q}_M, \mathbb{W}}$  and analyzes it in order to discover whether it is satisfied.

The assignment in line 4 computes the test set  $\text{TEST}_{\langle \mu(e), \mathfrak{q} \rangle}$  which terminates by Definition 13: actually, it is just an application of the simulation algorithm [24] to a finite set of (marked) terms, i.e., the Web site  $\mathbb{W}$ .

The emptiness test in lines 5–7 trivially terminates as well as the *case* statement in lines 8–17, which includes some checks based on the simulation algorithm [24].

Now, let us prove the partial correctness of the algorithm: that is, if an error message regarding  $e$  (incompleteness symptom) is returned by the procedure, then  $\langle \mu(e), \mathfrak{q} \rangle$ , where  $\mathfrak{q} \in \{A, E\}$ , is not satisfied in  $\mathbb{W}$ .

If an error message regarding  $e$  (incompleteness symptom) is returned by the procedure, then it refers to a requirement  $\langle \mu(e), \mathfrak{q} \rangle$ . Let us consider the iteration of the loop in lines 3–18 which checks  $\langle \mu(e), \mathfrak{q} \rangle$ . Note that an error message can be risen in three cases. Whenever

1.  $\text{TEST}_{\langle \mu(e), \mathfrak{q} \rangle} = \emptyset$ .
2. There exists  $p \equiv (V, E, r, \text{label}) \in \text{TEST}_{\langle \mu(e), A \rangle}$  s.t.  $e \not\equiv p|_V$ , with  $v \in V$ .
3. For each  $p \equiv (V, E, r, \text{label}) \in \text{TEST}_{\langle \mu(e), E \rangle}$ ,  $e \not\equiv p|_V$ , with  $v \in V$ .

In case 1, requirement  $\langle \mu(\epsilon), \varrho \rangle$  is not satisfied in  $w$  directly by Definition 14 (see point 1) [resp., Definition 15 (see point 1)] when  $\varrho = A$  (resp.,  $\varrho = E$ ).

In case 2, requirement  $\langle \mu(\epsilon), \varrho \rangle$ , where  $\varrho = A$ , is not satisfied in  $w$  by Definition 14 (see point 2).

In case 3, requirement  $\langle \mu(\epsilon), \varrho \rangle$ , where  $\varrho = E$ , is not satisfied in  $w$  by Definition 15 (see point 2).

## References

- Abiteboul, S., Buneman, P., Suciu, D.: Data on the Web. From Relations to Semistructured Data and XML. Morgan Kaufmann (2000)
- Alpuente, M., Ballis, D., Falaschi, M.: A rewriting-based framework for Web sites verification. In: Proceedings of 1st International Workshop on Ruled-Based Programming (RULE'04), vol. 124(1). ENTCS, Elsevier (2004)
- Alpuente, M., Ballis, D., Falaschi, M.: VERDI: an automated tool for Web sites Verification. In: Alferes, J.J., Leite, J. (eds) Proceedings of the 9th European Conference on Logics in Artificial Intelligence (JELIA'04), vol. 3229 of Lecture Notes in Computer Science, pp. 726–729. Springer, Berlin Heidelberg New York (2004)
- Baader, F., Nipkow, T.: Term Rewriting and All That. Cambridge University Press, Cambridge (1998)
- Ballis, D.: Rule-based Software Verification and Correction. PhD thesis, University of Udine and Technical University of Valencia (2005)
- Ballis, D., García Vivó, J.: A rule-based system for Web site verification. In: Proceedings of 1st International Workshop on Automated Specification and Verification of Web Sites (WWV'05). ENTCS, Elsevier, 2005. To appear.
- Baxter, I.D., Ricca, F., Tonella, P.: Web application transformations based on rewrite rules. *Inform. Softw. Technol.* **44** (13), (2002)
- Bertino, E., Mesiti, M., Guerrin, G.: A matching algorithm for measuring the structural similarity between an XML document and a DTD and its applications. *Inform. Syst.* **29**(1), 23–46 (2004)
- Bezem, M.: TeReSe, Term Rewriting Systems, chapter Mathematical background (Appendix A). Cambridge University Press, Cambridge (2003)
- Bry, F., Schaffert, S.: Towards a declarative query and transformation language for XML and semistructured data: simulation unification. In: Proceedings of the International Conference on Logic Programming (ICLP'02) vol. 2401 of LNCS. Springer Berlin Heidelberg New York (2002)
- Bry, F., Schaffert, S.: The XML query language xcerpt: design principles, examples, and semantics. Technical report, Available at: <http://www.xcerpt.org> (2002)
- Buneman, P., Davidson, S.B., Hillebrand, G.G., Suciu, D.: A query language and optimization techniques for unstructured data. In: Proceedings of ACM SIGMOD International Conference on Management of Data (ICMD'96) (1996)
- Capra, L., Emmerich, W., Finkelstein, A., Nentwich, C.: XLINKIT: a consistency checking and smart link generation service. *ACM Trans. Internet Technol.* **2** (2), 151–185 (2002)
- Cortesi, A., Dovier, A., Quintarelli, E., Tanca, L.: Operational and abstract semantics of a graphical query language. *Theor. Comput. Sci.* **275**, 521–560 (2002)
- Dershowitz, N., Plaisted, D.: Rewriting. *Handbook Automated Reasoning* **1**, 535–610 (2001)
- Despeyroux, T., Trousse, B.: Semantic verification of Web sites using natural semantics. In: Proceedings of 6th Conference on Content-Based Multimedia Information Access (RIAO'00) (2000)
- Di Sciascio, E., Donini, F.M., Mongiello, M., Piscitelli, G.: Web applications design and maintenance using symbolic model checking. In: Proceedings 7th European Conference on Software Maintenance and Reengineering, pp. 63. IEEE Computer Society (2003)
- Easterbrook, S.M., Nuseibeh, B., Russo, A.: Leveraging inconsistency in software development. *IEEE Comp.* **33** (4), 24–29 (2000)
- Ellmer, E., Emmerich, W., Finkelstein, A., Nentwich, C.: Flexible consistency checking. *ACM Trans. Softw. Eng.* **12**(1), 28–63 (2003)
- Fan, W., Libkin, L.: On XML integrity constraints in the presence of DTDs. *J. ACM* **49**(3), 368–406 (2002)
- Fernandez, M., Florescu, D., Levy, A., Suciu, D.: Verifying integrity constraints on Web sites. In: Proceedings of Sixteenth International Joint Conference on Artificial Intelligence (IJCAI'99) vol. 2 pp. 614–619. Morgan Kaufmann (1999)
- Fernandez, M.F., Suciu, D.: Optimizing regular path expressions using graph schemas. In: Proceedings of International Conference on Data Engineering (ICDE'98), pp. 14–23 (1998)
- M. Hanus (ed.). Curry: an integrated functional logic language. Available at: <http://www-i2.informatik.rwthachen.de/~hanus/curry> (1999)
- Henzinger, M.R., Henzinger, T.A., Kopke, P.W.: Computing simulations on finite and infinite graphs. In: IEEE Sympo Foundations Comput. Sci. 453–462 (1995)
- Hopcroft, J.E., Ullman, J.D.: Introduction to Automata Theory, Languages, and Computation. Addison-Wesley (1979)
- Hosoya, H., Pierce, B.: Regular expressions pattern matching for XML. In: Proceedings of 25th ACM SIGPLAN-SIGACT International Symposium POPL, pp. 67–80. ACM (2001)
- Hussmann, H.: Unification in conditional-equational Theories. In: Proceedings of ALP'88, pp. 31–40. Springer LNCS 343 (1988)
- Imagware, Inc. Doctor HTML: quality assessment for the Web. Available at: <http://www.doctor-html.com/RxHTML/>.
- Kirchner, C., Qian, Z., Singh, P.K., Stuber, J.: Xemantics: a Rewriting Calculus-Based Semantics of XSLT. Rapport de recherche A01-R-386, LORIA (2001)
- Klop, J.W.: Term rewriting systems. In: Abramsky, S., Gabbay, D., Maibaum, T., (eds.) Handbook of Logic in Computer Science vol. I, pp. 1–112. Oxford University Press, Oxford (1992)
- Moreno-Navarro, J.J., Rodriguez-Artalejo: BABEL: a functional and logic programming language based on a constructor discipline and narrowing. In: Grabowski, I., Lescanne, P., Wechler, W., (eds.) Proceedings of the International Conference on Algebraic and Logic Programming, pp. 223–232 Springer LNCS 343 (1988)
- Nentwich, C., Emmerich, W., Finkelstein, A.: Consistency management with repair actions. In: Proceedings of the 25th International Conference on Software Engineering (ICSE'03). IEEE Computer Society (2003)
- Nesbit, S.: HTML Tidy: keeping it clean Available at: <http://www.webreview.com/2000/06-16/webauthors/06-16-00-3.shtml> (2000)
- The Open Group. Unix Regular Expressions. Available at: <http://www.opengroup.org/onlinepubs/7908799/xbd/re.html>.

35. Academia Sinica Computing Centre. The schematron: an XML structure validation language using pattern in trees. Available at: <http://xml.ascc.net/resource/schematron/schematron.html>
36. World Wide Web Consortium (W3C). Extensible Markup Language (XML) 1.0, second edition Available at: <http://www.w3.org> (1999)
37. World Wide Web Consortium (W3C). XML path language (XPath) Available at: <http://www.w3.org> (1999)
38. World Wide Web Consortium (W3C). Extensible HyperText Markup Language (XHTML) Available at: <http://www.w3.org> (2000)