

A Rewriting-based Framework for Web Sites Verification^{*}

M. Alpuente¹, D. Ballis², and M. Falaschi²

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

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

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 obtains the requirements not fulfilled by the Web site, and helps to repair the errors by finding out incomplete information and/or missing pages. 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 Web verification system VERDI 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 quality 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 repair errors by finding out missing pages and/or incomplete information, such as the data or the links available in a particular page.

Although the management of Web sites has received significant attention in recent years [3, 8, 9], few works address the semantic verification of Web sites. In [9], a declarative verification algorithm is proposed 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. [8] proposes a methodology which

^{*} This work has been partially supported by CICYT under grant TIC2001-2705-C03-01, by Agencia Valenciana de Ciencia y Tecnología, under grant GR03/025

consists of using inference rules and axioms to define some semantic constraints concerning the Web site contents. Then, a verification technique is proposed which is based on compiling the specification into Prolog code. 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 and repairing techniques. As far as we know, rewriting-based techniques have not been explored in this context to date.

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 enforce that some information is available at a given Web page, some links between pages do exist 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 incomplete and/or missing Web pages. Moreover, by analyzing the requirements not fulfilled by the Web site, we are also able to find out the missing information which is needed to repair the Web site. Since reasoning on the Web calls for formal methods specifically fitting the Web context, we develop a novel, rewriting-based technique called *partial rewriting*, in which the traditional pattern matching mechanism is replaced with tree *simulation* [11] in order to provide a suitable mechanism for recognizing patterns inside semistructured documents. The notion of simulation has been already used before for dealing with semistructured data in a number of query and transformation languages [3, 5, 10, 6]. The reason is twofold: on the one hand, it provides a powerful method to extract information from semistructured data; on the other hand, efficient algorithms exist for computing simulations [11]. To assess the feasibility and efficiency of our approach, we have implemented the prototype system VERDI (VERification and Rewriting for Debugging Internet sites), which is based on the verification methodology that we propose and is publicly available online.

Plan of the paper. Section 2 summarizes some preliminary definitions and notations. In Section 3, we formalize a quite obvious method for translating HTML/XML documents into Herbrand terms. Section 4 illustrates the notion of *page simulation*, whereas Section 5 is devoted to formalize the specification language as well as the *partial rewriting* mechanism. In Section 6, we introduce our verification technique, which is formalized as a fixpoint computation. First, the set of requirements to be fulfilled by the Web site W are computed as the fixpoint of a suitable operator associated with the Web site specification I . Then, by using simulation we select those requirements which are not satisfied by W and the corresponding incomplete/missing Web pages which are the source for the errors. The requirements which are not satisfied also allow us to ascertain the missing information which is needed to repair the Web site. Some notes

regarding the implementation of the system VERDI are given in Section 7. Section 8 concludes. Proofs of the main results are included in Appendix 8.

2 Preliminaries

We call *alphabet* a finite set of symbols. 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 Σ denotes a set of function symbols, or *signature*. We consider varyadic signatures as in [7]. $\tau(\Sigma, \mathcal{V})$ and $\tau(\Sigma)$ denote the *non-ground term algebra* and the *term algebra* built on $\Sigma \cup \mathcal{V}$ and Σ , respectively. $\tau(\Sigma)$ is usually called the Herbrand universe over Σ . A term t is *linear*, if no variable appears more than once in t .

Terms are viewed as labelled trees in the following way: a term in $\tau(\Sigma)$ is a tree $(V, E, r, label)$, 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 $label(v) \in (\Sigma \cup \mathcal{V})$, for each $v \in V$. Let us see a small example.

Example 1. Consider the term $t \equiv (f(g(a), X))$ in $\tau(\{f, g, a\}, \{X\})$. Term t can be represented by the structure $(V, E, r, label)$, 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: $label(v_0) = f$, $label(v_1) = g$, $label(v_2) = X$, $label(v_3) = a$.

Given two vertices $v, v' \in V$ of a term $t \equiv (V, E, r, label)$, 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 . 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 $depth(t, v)$. A *substitution* $\sigma \equiv \{X_1/t_1, X_2/t_2, \dots\}$ is a mapping from the set of variables \mathcal{V} into the set of terms $\tau(\Sigma, \mathcal{V})$. By $Var(t)$ we denote the set of variables occurring in term t .

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

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

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

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

Given a term $t \equiv (V, E, r, label)$ and a valid marking μ for t , by slightly abusing notation we recursively define a *marked* term $\mu(t)$ as follows:

$$\mu(t) = \begin{cases} \underline{X} & t \equiv (\{v\}, \emptyset, v, label) \wedge label(v) = X \in \mathcal{V} \\ & \wedge \mu(v) = yes \\ X & t \equiv (\{v\}, \emptyset, v, label) \wedge label(v) = X \in \mathcal{V} \\ & \wedge \mu(v) = no \\ \underline{f}(\mu(t_1), \dots, \mu(t_n)) & t \equiv (V, E, r, label) \equiv f(t_1, \dots, t_n) \wedge \mu(r) = yes \\ \underline{\bar{f}}(\mu(t_1), \dots, \mu(t_n)) & t \equiv (V, E, r, label) \equiv f(t_1, \dots, t_n) \wedge \mu(r) = no \end{cases}$$

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

Example 2. Consider again term $t \equiv (f(g(a), X))$ of Example 1. Let μ_1 be a marking for t defined as $\mu_1(v_0) = \mu_1(v_2) = \mu_1(v_3) = yes$, $\mu_1(v_1) = no$. Additionally, let μ_2 be a marking for t such that $\mu_2(v_0) = \mu_2(v_1) = yes$, $\mu_2(v_2) = \mu_2(v_3) = no$. Note that μ_1 is not a valid marking for t as the marked part of t is not a term in $\tau(\{f, g, a\}, \{X\})$, whereas μ_2 is valid for t and $\mu_2(t) = \underline{f}(g(a), X)$ is a marked term.

3 Denotation of Web Sites

In this paper, a *Web page* is either a XML[14] or a HTML[13] document, and a *Web site* is a finite collection of Web pages. In the sequel, we provide a formalization of these concepts by means of semistructured expressions, which can be seen as an abstract syntax which generalizes the two markup languages XML and HTML. Then, we show how semistructured expressions can be translated into ordinary terms of a given term algebra in such a way that Web sites are represented as finite sets of (ground) terms.

Semistructured Expressions. XML/HTML documents consist of nested structured data, which can be defined recursively. Abstracting from XML and HTML, we give a formal definition of semistructured expressions which are suitable for representing structured documents written in one of these two languages.

Let us consider two alphabets T and $\mathcal{T}ag$. We denote the set T^* by $\mathcal{T}ext$. An object $\mathbf{t} \in \mathcal{T}ag$ is called *tag* element, while an element $\mathbf{w} \in \mathcal{T}ext$ is called *text* element. A *semistructured expression* \mathbf{e} over $\mathcal{T}ext$ and $\mathcal{T}ag$ sets can be specified by the following syntax³

$$\begin{aligned} \mathbf{e} & := \langle \mathbf{t} \rangle \mathbf{elist} \langle / \mathbf{t} \rangle \mid \mathbf{w} \quad \forall \mathbf{w} \in \mathcal{T}ext, \mathbf{t} \in \mathcal{T}ag \\ \mathbf{elist} & := \mathbf{eelist} \mid \epsilon \end{aligned}$$

We denote the set of all the semistructured expressions over $\mathcal{T}ext$ and $\mathcal{T}ag$ by $\mathcal{S}(\mathcal{T}ext, \mathcal{T}ag)$. Note that $\mathcal{T}ext \subseteq \mathcal{S}(\mathcal{T}ext, \mathcal{T}ag)$.

³ Note that symbol ϵ in the syntax given for semistructured expressions denotes the empty string and must not be confused with the empty marking ε .

Example 3. The following object is a semistructured expression.

```

<members>
  <member>
    <name> mario </name>
    <surname> rossi </surname>
    <status> professor </status>
  </member>
  <member>
    <name> franca </name>
    <surname> bianchi </surname>
    <status> technician </status>
  </member>
  <member>
    <name> giulio </name>
    <surname> verdi </surname>
    <status> student </status>
  </member>
</members>

```

Roughly speaking, a semistructured expression is either a raw or a structured piece of text, where the structure is provided by tags. Consequently, tags allow us to mark up some textual content, which may contain an arbitrary amount of further well-bracketed markup. Informally, the more tags we add, the more the text is structured, and in some sense its “formal organization” will also increase. Note that we have not explicitly dealt with XML/HTML attributes, as they can be seen as common tagged elements and thus modeled as semistructured expressions. On the other hand, without loss of generality, other XML/HTML features such as namespaces, DTDs and/or schemas, that are not relevant to this work are not conveyed by our notion of semistructured expression.

In the literature, slightly different formalisms have been introduced for modeling XML/HTML documents, e.g. in [1] semistructured expressions are directed graphs which can deal with crossing references. Nevertheless, we prefer the hierarchical representation which does not cause any serious restriction in many practical contexts while greatly simplifies our methodology.

Term representation. Semistructured expressions are provided with a tree-like structure, therefore they can be conveniently translated into terms by applying the following straightforward transformation.

Definition 1. *Let e be a semistructured expression over \mathcal{Text} and \mathcal{Tag} . Then, e is represented by a term of the Herbrand universe $\tau(\mathcal{Text} \cup \mathcal{Tag})$ by the translation $s_to_t: \mathcal{S}(\mathcal{Text}, \mathcal{Tag}) \rightarrow \tau(\mathcal{Text} \cup \mathcal{Tag})$ defined as follows:*

$$s_to_t(e) = \begin{cases} w & \text{if } e \equiv w \in \mathcal{Text} \\ \mathfrak{t}(s_to_t(e_1), \dots, s_to_t(e_n)) & \text{if } e \equiv \langle \mathfrak{t} \rangle e_1 \dots e_n \langle / \mathfrak{t} \rangle \end{cases}$$

Example 4. Consider again semistructured expression of Example 3. Then, the term p computed by function *s-to-t* for that semistructured expression is

```
members(
  member(name(mario), surname(rossi), status(professor)),
  member(name(franca), surname(bianchi), status(technician)),
  member(name(giulio), surname(verdi), status(Student))
)
```

To summarize, a Web page, which is coded as a HTML/XML document, can be represented as a semistructured expression, which is then easily translated into a corresponding term of a suitable term algebra. Therefore, in the remaining of this work, a Web page is modeled by a term in $\tau(\text{Text} \cup \text{Tag})$. Besides, a *marked* Web page is defined as $\mu(p)$, where $p \in \tau(\text{Text} \cup \text{Tag})$ and μ 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 μ is a valid marking for s . In our methodology, (marked) Web page templates are used for specifying properties on Web sites as we will see later on.

4 Page Simulations

In this section, we formalize a notion of *page simulation* for Web pages which 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 labelled 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 5. 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 Figure 1, 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 nodes of 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 nodes of 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 .

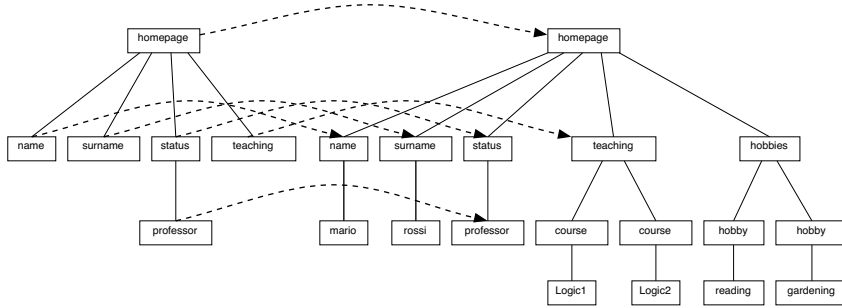


Fig. 1. Page simulation between p_1 and p_2 .

Simulations have been used in a number of works dealing with querying and transformation of semistructured data. For instance, [1, 10] propose some techniques based on simulation for analyzing semistructured data w.r.t. a given schema. The language Xcerpt [4, 3] is a (logic) query language for XML and semistructured documents which implements a sort of unification by exploiting the notion of graph simulation.

Basically, the reason why simulations are successfully employed in the implementation of these kinds of manipulation and querying methods is twofold. Firstly, it is a simple and powerful technique to extract and recognize the partial structure of a document; secondly, there are several efficient algorithms to compute (graph and tree) simulations (see [11]).

In the following, we provide our notion of simulation which is a slight adaptation of the one given in [3] 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 2. 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 $\tau(\mathcal{Text} \cup \mathcal{Tag}, \mathcal{V})$. The label relation $\sim \subseteq V_1 \times V_2$ is defined as follows:

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

Definition 3. Let $s_1 \equiv (r_1, V_1, E_1, label_1)$, $s_2 \equiv (r_2, V_2, E_2, label_2)$ be two Web page templates in $\tau(\mathcal{Text} \cup \mathcal{Tag}, \mathcal{V})$ 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 $\pi(\mathbf{S}) = \{v_2 \mid (v_1, v_2) \in \mathbf{S}\}$.

Roughly speaking, Definition 3 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 3 formalizes the similarity w.r.t labelings, that is, any pair of nodes (v, v') in a page simulation \mathbf{S} of \mathbf{s}_1 in \mathbf{s}_2 have the same label, otherwise node v must be labelled by a variable, which somehow means that the label of v can be seen as a generalization of any concrete label of v' . Finally, Condition (1) and Condition (3) provide a relation between the tree structure of \mathbf{s}_1 and the tree structure of \mathbf{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.

In the following, we define a subclass of simulations called minimal simulations.

Definition 4. Let $\mathbf{s}_1 \equiv (V_1, E_1, r_1, label_1)$, $\mathbf{s}_2 \equiv (V_2, E_2, r_2, label_2)$ be two Web page templates in $\tau(\text{Text} \cup \text{Tag}, \mathcal{V})$. A page simulation \mathbf{S} of \mathbf{s}_1 in \mathbf{s}_2 w.r.t. \sim is minimal if there are no page simulations \mathbf{S}' of \mathbf{s}_1 in \mathbf{s}_2 w.r.t. \sim such that $\mathbf{S}' \subseteq \mathbf{S}$.

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

Example 6. Let us consider the following Web page templates \mathbf{s}_1 and \mathbf{s}_2 : $\text{hobbies}(\text{hobby}(X))$, $\text{hobbies}(\text{hobby}(\text{reading}), \text{hobby}(\text{gardening}))$. In Figure 2(a), the dashed arrows represent a non-minimal simulation of \mathbf{s}_1 in \mathbf{s}_2 , while in Figures 2(b) and 2(c) two minimal simulations of \mathbf{s}_1 in \mathbf{s}_2 are depicted. Note that the last two simulations are mappings.

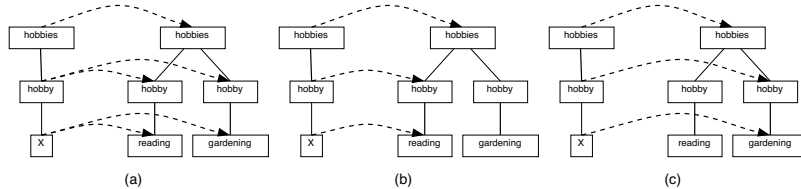


Fig. 2. non-minimal and minimal simulations

Lemma 1. Let $\mathbf{s}_1 \equiv (V_1, E_1, r_1, label_1)$, $\mathbf{s}_2 \equiv (V_2, E_2, r_2, label_2)$ be two Web page templates in $\tau(\text{Text} \cup \text{Tag}, \mathcal{V})$. A minimal page simulation \mathbf{S} of \mathbf{s}_1 in \mathbf{s}_2 w.r.t. \sim is a total mapping $\mathbf{S} : V_1 \rightarrow V_2$.

Minimal simulations do not guarantee that the tree structure of a given Web page template can be recognized inside another template. For this purpose, we need to furtherly restrict our class of simulations. Let us see an example.

Example 7. Consider Web page templates $\mathbf{s}_1 \equiv \mathbf{f}(X, Y)$ and $\mathbf{s}_2 \equiv \mathbf{f}(\mathbf{a})$. Note that there exists a minimal page simulation of \mathbf{s}_1 in \mathbf{s}_2 w.r.t. \sim (see Figure 3), but the tree structure of \mathbf{s}_1 cannot be recognized as part of \mathbf{s}_2 , e.g. the vertex with label \mathbf{f} in \mathbf{s}_1 has two outgoing edges, while the corresponding vertex in \mathbf{s}_2 has only one.

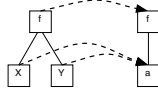


Fig. 3. minimal non-injective simulation

To solve the problem presented in Example 7, we simply restrict ourselves to consider minimal *injective* page simulations, which provide a one-to-one correspondence among edges of the two considered Web page templates.

It is not difficult to demonstrate that minimal injective simulations are particular instances of Kruskal’s *embeddings* [2] w.r.t. the relation \sim . In other words, a minimal injective page simulation of \mathbf{s}_1 in \mathbf{s}_2 w.r.t. \sim exists iff \mathbf{s}_1 is embedded into \mathbf{s}_2 w.r.t. \sim , i.e., we are able to find out the structure and the labeling of \mathbf{s}_1 inside \mathbf{s}_2 . Note that the minimal simulation of \mathbf{s}_1 in \mathbf{s}_2 depicted in Figure 3 is not injective and thus no embedding of \mathbf{s}_1 into \mathbf{s}_2 exists. Instead, Figures 2(b) and 2(c) illustrate two minimal injective simulations, that is, two embeddings between Web page templates.

5 Web specification language

In this section, 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 specification is a finite set of rules, where the terms in the left-hand side and in the right-hand side of each rule represent (eventually marked) Web page templates. The operational mechanism we consider is based on a novel rewriting-based mechanism, which is able to extract partial structure from a term, and rewrite it by using page simulation. Formally, Web site specifications are as follows.

Definition 5. A rule is a pair of terms $\mathbf{l} \rightarrow \mu(\mathbf{r})$ such that $\mathbf{l}, \mathbf{r} \in \tau(\text{Text} \cup \text{Tag}, \mathcal{V})$, \mathbf{l} is linear, $\text{Var}(\mathbf{r}) \subseteq \text{Var}(\mathbf{l})$ and μ is a valid marking for \mathbf{r} . A Web site specification \mathbf{I} is a finite set of rules $\{\mathbf{l}_1 \rightarrow \mu_1(\mathbf{r}_1), \dots, \mathbf{l}_n \rightarrow \mu_n(\mathbf{r}_n)\}$.

Given a Web specification \mathbf{I} , we denote the set of all left-hand sides (right-hand sides disregarding markings) of rules in \mathbf{I} by $\text{Lhs}_{\mathbf{I}}$ ($\text{Rhs}_{\mathbf{I}}$, respectively). In symbols, $\text{Lhs}_{\mathbf{I}} = \{\mathbf{l} \mid \mathbf{l} \rightarrow \mu(\mathbf{r}) \in \mathbf{I}\}$ and $\text{Rhs}_{\mathbf{I}} = \{\mathbf{r} \mid \mathbf{l} \rightarrow \mu(\mathbf{r}) \in \mathbf{I}\}$.

Rules of a Web specification formalize conditions to be fulfilled by a given Web site. Intuitively, the interpretation of a rule $\mathbf{l} \rightarrow \mu(\mathbf{r})$ w.r.t. a Web site \mathbf{W}

is as follows: if (an instance of) l is recognized in W , also (an instance of) r must be recognized in a subset of W , which is determined by computing the sets of all Web pages which embed (an instance of) the marked part of r .

Roughly speaking, markings in the right-hand sides of the rules allow us to find sets of Web pages, which might be incomplete or missing. Then, real buggy pages are detected inside these sets.

In the following example, we informally illustrate the definition of Web specification. Marks are introduced by the user to help locating errors. We do not take care of marks by now but postpone the formal handling of marking information and the description of the verification framework to Section 6.

Example 8. Consider the following Web specification, which models some required properties of a research group Web site containing information about group members affiliation, scientific publications and personal data.

$$\begin{aligned} \underline{\text{member}}(\text{name}(X), \text{surname}(Y)) &\rightarrow \underline{\text{hpage}}(\text{name}(X), \underline{\text{surname}}(Y), \text{status}) \\ \text{hpage}(\text{status}(\text{professor})) &\rightarrow \underline{\text{hpage}}(\underline{\text{status}}(\text{professor}), \text{teaching}) \\ \text{pubs}(\text{pub}(\text{name}(X), \text{surname}(Y))) &\rightarrow \underline{\text{member}}(\text{name}(X), \text{surname}(Y)) \end{aligned}$$

First rule formalizes the following property: if there is a Web page containing a member list, then for each member, a home page exists containing (at least) the name, the surname and the status of this member. Second rule states that whenever a home page of a professor is recognized, then that page must also include some teaching information. Finally, the 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 order to mechanize the intended semantics of Web specification rules, we have to 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.

Definition 6. Let $\mathbf{s}_1 \equiv (V_1, E_1, r_1, \text{label}_1)$, $\mathbf{s}_2 \equiv (V_2, E_2, r_2, \text{label}_2) \in \tau(\text{Text} \cup \text{Tag}, \mathcal{V})$. We say that \mathbf{s}_2 partially matches \mathbf{s}_1 via substitution σ iff

1. there exists a minimal injective simulation \mathbf{S} of \mathbf{s}_1 in \mathbf{s}_2 w.r.t. \sim ;
2. for each $(v, v') \in \mathbf{S}$ such that $\text{label}(v) = X \in \mathcal{V}$, $\sigma(X) = (\mathbf{s}_2|_{v'})$.

In Definition 6, we consider only minimal injective simulations between Web page templates \mathbf{s}_1 and \mathbf{s}_2 , since this trivially ensures the existence of a substitution σ such that there exists a simulation of $\mathbf{s}_1\sigma$ in \mathbf{s}_2 w.r.t. \sim ; in other words, $\mathbf{s}_1\sigma$ is embedded into \mathbf{s}_2 .

Example 9. Consider again Example 6. We have that \mathbf{s}_2 partially matches \mathbf{s}_1 via $\{X/\text{reading}\}$ (see Figure 2(b)) and \mathbf{s}_2 partially matches \mathbf{s}_1 via $\{X/\text{gardening}\}$ (see Figure 2(c)). Note that performing partial matching by the non-minimal simulation of Figure 2(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.

Definition 7. Let $\mathbf{s} \equiv (V, E, r, \text{label})$, $\mathbf{t} \in \tau(\text{Text} \cup \text{Tag}, \mathcal{V})$. Let μ_1 and μ_2 be two valid markings for \mathbf{s} and \mathbf{t} , respectively. Then, $\mu_1(\mathbf{s})$ partially rewrites to $\mu_2(\mathbf{t})$ via rule $r \equiv \mathbf{l} \rightarrow \mu(\mathbf{r})$ and substitution σ (in symbols, $\mu_1(\mathbf{s}) \rightarrow_r^\sigma \mu_2(\mathbf{t})$) iff there exists $v \in V$ such that

1. \mathbf{l} partially matches $\mathbf{s}|_v$ via σ ;
2. $\mathbf{t} = \mathbf{r}\sigma$.
3. Let $\mathbf{r} \equiv (V_{\mathbf{r}}, E_{\mathbf{r}}, r, \text{label}_{\mathbf{r}})$ and $\mathbf{r}\sigma \equiv (V_{\mathbf{r}\sigma}, E_{\mathbf{r}\sigma}, r, \text{label}_{\mathbf{r}\sigma})$. For each $v \in V_{\mathbf{r}\sigma}$,

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

When rule r and substitution σ are understood, we simply write $\mu_1(\mathbf{s}) \rightarrow \mu_2(\mathbf{t})$.

It is worth noting that we provide a notion of partial rewriting in which the context of the selected reducible expression $\mathbf{s}|_v$ of the Web page template which is rewritten is disregarded after the rewrite step (see point (2) of Definition 7). Roughly speaking, given a Web specification rule $\mathbf{l} \rightarrow \mu(\mathbf{r})$, partial rewriting allows us to extract a subpart of a given Web page (template) \mathbf{s} , which partially matches \mathbf{l} , and to replace \mathbf{s} by an instance of \mathbf{r} ; namely, $\mathbf{r}\sigma$ (see points (1) and (2) of Definition 7). Point (3) of Definition 7 establishes that rewritten templates inherit markings from the right-hand sides of the applied rules. More precisely,

- each vertex of $\mathbf{r}\sigma$, which is not affected by substitution σ , maintains the same marking of \mathbf{r} ;
- each vertex, which belongs to a subterm of $\mathbf{r}\sigma$ replacing a variable \underline{X} of \mathbf{r} , is marked *yes*;
- each vertex, which belongs to a subterm of $\mathbf{r}\sigma$ replacing a variable X of \mathbf{r} , is marked *no*.

Example 10. Consider the Web page \mathbf{p} of Example 4 and the first rule \mathbf{r}_1 of the Web specification of Example 8. Then, Web page template $\varepsilon(\mathbf{p})$ partially rewrites to the following three Web pages by applying \mathbf{r}_1 .

$$\begin{aligned} \varepsilon(\mathbf{p}) &\rightarrow_{\mathbf{r}_1} \underline{\text{hpage}}(\text{name}(\text{mario}), \underline{\text{surname}}(\text{rossi}), \text{status}) \\ \varepsilon(\mathbf{p}) &\rightarrow_{\mathbf{r}_1} \underline{\text{hpage}}(\text{name}(\text{franca}), \underline{\text{surname}}(\text{bianchi}), \text{status}) \\ \varepsilon(\mathbf{p}) &\rightarrow_{\mathbf{r}_1} \underline{\text{hpage}}(\text{name}(\text{giulio}), \underline{\text{surname}}(\text{verdi}), \text{status}) \end{aligned}$$

6 The verification framework

In this section, we show how simulation and partial rewriting can be applied to verify a given Web site \mathbf{W} w.r.t. a Web specification \mathbf{I} . Essentially, the main idea is to compute the set of all possible marked Web pages that can be derived from \mathbf{W} via \mathbf{I} by means of partial rewriting. These marked Web pages can be thought

of as requirements to be fulfilled by W . Then, we check whether the computed requirements are satisfied by W by using simulation and marking information. In summary, the method works in two steps, which are repeatedly applied as we described in the following.

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

6.1 Computing the set of requirements

Let us introduce the following operator.

Definition 8. Let \mathbf{T} be a set of marked Web page templates and \mathbf{I} be a Web specification. Then,

$$\mathbf{R}_{\mathbf{I}}(\mathbf{T}) = \mathbf{T} \cup \{\mu_2(\mathbf{s}_2) \mid \exists \mu_1(\mathbf{s}_1) \in \mathbf{T}, r \equiv 1 \rightarrow \mu(\mathbf{r}) \in \mathbf{I} \text{ s.t. } \mu_1(\mathbf{s}_1) \rightarrow_r \mu_2(\mathbf{s}_2)\}$$

Roughly speaking, the operator in Definition 8 computes all marked templates which result from partial rewriting the Web page templates of \mathbf{T} by using the Web specification \mathbf{I} , and returns the union of the resulting set and \mathbf{T} . By repeatedly applying this operator, it is possible to compute all marked Web pages that can be derived from an initial Web site after an arbitrary number of partially rewriting steps. For this purpose, we formalize the *ordinal powers* of the operator $\mathbf{R}_{\mathbf{I}}$ w.r.t. a Web site W as follows: $\mathbf{R}_{\mathbf{I}} \uparrow^W 0 = W$, $\mathbf{R}_{\mathbf{I}} \uparrow^W n = \mathbf{R}_{\mathbf{I}}(\mathbf{R}_{\mathbf{I}} \uparrow^W (n-1))$, $n > 0$.

It is immediate to demonstrate that the operator $\mathbf{R}_{\mathbf{I}}$ is continuous on the lattice consisting of the powerset of $\tau(\text{Text} \cup \text{Tag}, \mathcal{V})$ ordered by set inclusion. This ensures that a least fixpoint of $\mathbf{R}_{\mathbf{I}}$ exists and can be reached after ω applications of $\mathbf{R}_{\mathbf{I}}$, that is, $\mathbf{R}_{\mathbf{I}} \uparrow^W \omega$ where ω is the first infinite ordinal. Moreover, the least fixpoint of $\mathbf{R}_{\mathbf{I}}$ contains all the marked Web pages derivable from Web pages in W via \mathbf{I} .

Now, recalling the interpretation of the rules of the Web site specification given in Section 5, Web pages derived by the application of a Web specification must be recognized as (part of) some Web page in the Web site. Therefore, those Web pages in the least fixpoint of $\mathbf{R}_{\mathbf{I}}$ which are not in W can be intended as requirements to be fulfilled by W . Thus, we define the *set of requirements* for W w.r.t. \mathbf{I} as

$$\text{Req}_{\mathbf{I},W} = \text{lfp}(\mathbf{R}_{\mathbf{I}}) \setminus W$$

where $\text{lfp}(\mathbf{R}_{\mathbf{I}})$ is the least fixpoint of the operator $\mathbf{R}_{\mathbf{I}}$.

Clearly, the fixpoint of $\mathbf{R}_{\mathbf{I}}$ (and hence $\text{Req}_{\mathbf{I},W}$) for an arbitrary Web specification might be infinite. Consider for instance the following example.

Example 11. Let $W \equiv \{\mathbf{h}(\mathbf{g}(0), \mathbf{f}(0))\}$ be a Web site and $\mathbf{I} \equiv \{\mathbf{h}(\mathbf{g}(\mathbf{X})) \rightarrow \mathbf{h}(\mathbf{g}(\mathbf{g}(\mathbf{X})))\}$ be a Web specification. Then,

$$\text{Req}_{\mathbf{I},W} = \{\mathbf{h}(\mathbf{g}(\mathbf{g}(0))), \mathbf{h}(\mathbf{g}(\mathbf{g}(\mathbf{g}(0))))\}, \dots\}$$

is an infinite set of requirements which is not computable.

Fortunately, the computation of the set of requirements is finite for some interesting classes of Web specifications. Trivially, non-recursive specifications allow to reach $\text{lfp}(\mathbf{R}_I)$ after a finite number of applications of \mathbf{R}_I , i.e., $\text{lfp}(\mathbf{R}_I) = \mathbf{R}_I \uparrow^w k$, $k \in \mathbb{N}$. 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 8 cannot be formulated by using a non-recursive specification.

In the following, we define a class of recursive Web specifications for which the set of requirements is finite. Basically, the idea is to consider those specifications for which the computation of the least fixpoint only generates Web pages whose size is bounded.

The following definition formalizes the considered class of Web site specifications.

Definition 9. *A Web specification I is bounded iff, for each $\mathbf{l} \equiv (V_1, E_1, r_1, \text{label}_1) \in \mathbf{Lhs}_I$, $\mathbf{r} \equiv (V_2, E_2, r_2, \text{label}_2) \in \mathbf{Rhs}_I$ and each minimal injective simulation \mathbf{S} of \mathbf{l} in $\mathbf{r}|_v$ w.r.t. \sim , $v \in V_2$, the following property holds*

$$\text{if } v_2 \in \pi(\mathbf{S}) \text{ and } \text{label}_2(v_2) \in \text{Var}(\mathbf{r}|_v), \text{ then for all } v_1 \in V_1 \text{ s.t. } \text{label}_1(v_1) \in \text{Var}(\mathbf{l}), \\ \text{depth}(\mathbf{r}|_v, v_2) = \text{depth}(\mathbf{l}, v_1).$$

Roughly speaking, Definition 9 states that, whenever a left-hand side \mathbf{l} of a rule is simulated by (a subterm of) the right-hand side \mathbf{r} of a (possibly different) rule, then no variables in the substructure of \mathbf{r} which is recognized by simulation must be located at positions which are deeper than all the positions of the variables in \mathbf{l} .

Example 12. Consider again the specification I in Example 11. The left-hand side of the rule $\mathbf{h}(\mathbf{g}(\mathbf{X})) \rightarrow \mathbf{h}(\mathbf{g}(\mathbf{g}(\mathbf{X})))$ is simulated by its own right-hand side. Moreover, variable \mathbf{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, I is not bounded.

Now, take into account specification

$$I' \equiv \{\mathbf{m}(\mathbf{n}(\mathbf{X})) \rightarrow \mathbf{h}(\mathbf{n}(\mathbf{X}), \mathbf{s}(\mathbf{s}(\mathbf{X}))), \mathbf{h}(\mathbf{n}(\mathbf{X})) \rightarrow \mathbf{m}(\mathbf{n}(\mathbf{X}), \mathbf{t})\}.$$

Then, $\mathbf{m}(\mathbf{n}(\mathbf{X}))$ is simulated by $\mathbf{m}(\mathbf{n}(\mathbf{X}), \mathbf{t})$ and $\mathbf{h}(\mathbf{n}(\mathbf{X}))$ is simulated by $\mathbf{h}(\mathbf{n}(\mathbf{X}), \mathbf{s}(\mathbf{s}(\mathbf{X})))$. 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 I' is bounded.

For bounded Web specifications, the least fixpoint of the operator \mathbf{R}_I is finite as stated by the next proposition. This provides an effective method for computing the set of requirements $\text{Req}_{I, \mathbf{W}}$.

Proposition 1. *Let I be a bounded Web specification and \mathbf{W} be a Web site. Then, there exists $k \in \mathbb{N}$ such that $\text{lfp}(\mathbf{R}_I) = \mathbf{R}_I \uparrow^w k$.*

Example 13. Consider the bounded Web specification I of Example 8 and the following Web site \mathbf{W} :

```

W = { members(member(name(mario), surname(rossi), status(professor)),
               member(name(franca), surname(bianchi), status(technician)),
               member(name(giulio), surname(verdi), status(student))),
       hpage(name(mario), surname(rossi), phone(3333), status(professor),
             hobbies(hobby(reading), hobby(gardening))),
       hpage(name(franca), surname(bianchi), status(technician), phone(5555)),
       hpage(name(anna), surname(gialli), status(professor), phone(4444),
             teaching(course(algebra))),
       pubs(pub(name(mario), surname(rossi), title(blahblah1), year(2003)),
            pub(name(anna), surname(gialli), title(blahblah2), year(2002))) }

```

Then, the set of computed requirements $\text{Req}_{I,W}$ is

```

{ hpage(name(mario), surname(rossi), status),
  hpage(name(franca), surname(bianchi), status),
  hpage(name(giulio), surname(verdi), status),
  hpage(status(professor), teaching),
  member(name(mario), surname(rossi)),
  member(name(anna), surname(gialli)),
  hpage(name(anna), surname(gialli), status) }

```

6.2 Checking requirements in Web sites

As we have seen in Section 4, simulation allows us to identify the structure of a given Web page (eventually, a template) into another. By taking advantage of this fact, we can develop a methodology, which is able to discover incompleteness errors in a given Web site w.r.t. a Web specification. Basically, the idea is to verify the consistency of the Web site w.r.t. the set of requirements. To accomplish this task, we first use simulation for checking whether requirements are embedded into some Web page of the considered Web site and then exploit marking information in order to diagnose incompleteness errors in the Web site.

More precisely, our analysis allows us to discover two kinds of incompleteness errors: (1) Web pages which are missing in a Web site w.r.t. a given Web specification, (2) Web pages which are incomplete w.r.t. a given Web specification.

Let us first consider the former class of errors.

Definition 10. *Let W be a Web site, I be a bounded Web specification and $\text{Req}_{I,W}$ be the set of requirements for W w.r.t. I . Let $\mu(\mathbf{e}) \in \text{Req}_{I,W}$. The likely missed information set w.r.t. $\mu(\mathbf{e})$ is defined as*

$$LMIS_{\mu(\mathbf{e})} = \{p \equiv (V, E, r, \text{label}) \in W \mid \text{there is a minimal injective simulation of } \text{mark}(\mathbf{e}, \mu) \text{ in } \mathbf{p}|_v \text{ w.r.t. } \sim, \text{ with } v \in V\}.$$

Definition 10 computes a subset of the Web site containing (potentially incomplete) Web pages for each requirement. From the above definition and definition of simulation, it is rather easy to derive the following proposition.

Proposition 2. *Let W be a Web site, I be a bounded Web specification and $\text{Req}_{I,W}$ be the set of requirements for W w.r.t. I . Let $\mu(\mathbf{e}) \in \text{Req}_{I,W}$. If $LMIS_{\mu(\mathbf{e})} = \emptyset$,*

then for all $\mathbf{p} \equiv (V, E, r, \text{label}) \in \mathbb{W}$, and $v \in V$, there is no minimal injective simulation of \mathbf{e} in $\mathbf{p}|_v$.

Proposition 2 states that, whenever the likely missed information set is empty for a given requirement $\mu(\mathbf{e})$, $\mu(\mathbf{e})$ is not recognized in any Web page of the Web site. In other words, that requirement identifies a missing element in the Web site.

Definition 11. Let \mathbb{W} be a Web site, \mathbb{I} be a bounded Web specification and $\text{Req}_{\mathbb{I}, \mathbb{W}}$ be the set of requirements for \mathbb{W} w.r.t. \mathbb{I} . Let $\mu(\mathbf{e}) \in \text{Req}_{\mathbb{I}, \mathbb{W}}$ and $\mathbf{p} \in \mathbb{W}$. Then, $\mu(\mathbf{e})$ is missing in \mathbb{W} w.r.t. \mathbb{I} iff $\text{LMIS}_{\mu(\mathbf{e})} = \emptyset$.

Let us see an example for clarifying our definitions.

Example 14. Consider again the set of requirements $\text{Req}_{\mathbb{I}, \mathbb{W}}$ computed in Example 13. Then, $\mu(\mathbf{e}) \equiv (\text{hpage}(\text{name}(\text{giulio}), \text{surname}(\text{verdi}), \text{status}))$ is missing in \mathbb{W} w.r.t. \mathbb{I} , since $\text{LMIS}_{\mu(\mathbf{e})} = \emptyset$. Indeed, the requirement $\mu(\mathbf{e})$ identifies a “group member” home page which does not appear in the Web site \mathbb{W} .

Let us consider now incompleteness errors which refer to incomplete pages, that is, Web pages in which some piece of information is lacking (e.g. missing items).

Definition 12. Let \mathbb{W} be a Web site, \mathbb{I} be a bounded Web specification and $\text{Req}_{\mathbb{I}, \mathbb{W}}$ be the set of requirements for \mathbb{W} w.r.t. \mathbb{I} . Let $\mu(\mathbf{e}) \in \text{Req}_{\mathbb{I}, \mathbb{W}}$ and $\mathbf{p} \in \mathbb{W}$. Then, $\mathbf{p} \equiv (V, E, r, \text{label})$ is incomplete w.r.t. $\mu(\mathbf{e})$ iff

- $\mathbf{p} \in \text{LMIS}_{\mu(\mathbf{e})}$;
- there is a minimal injective simulation of $\text{mark}(\mathbf{e}, \mu)$ in $\mathbf{p}|_v$ w.r.t. \sim , with $v \in V$, s. t. there is no minimal injective simulation of \mathbf{e} in $\mathbf{p}|_v$ w.r.t. \sim .

In this case, we will call $\mu(\mathbf{e})$ incompleteness symptom for \mathbf{p} .

Example 15. Recall the set of requirements $\text{Req}_{\mathbb{I}, \mathbb{W}}$ computed in Example 13. Then, consider the requirement $\mu_1(\mathbf{e}_1) \equiv (\text{hpage}(\text{status}(\text{professor}), \text{teaching}))$, we have that

$$\text{LMIS}_{\mu_1(\mathbf{e}_1)} = \{(1) \text{hpage}(\text{name}(\text{mario}), \text{surname}(\text{rossi}), \text{phone}(3333), \text{status}(\text{professor}), \text{hobbies}(\text{hobby}(\text{reading}), \text{hobby}(\text{gardening}))), \\ (2) \text{hpage}(\text{name}(\text{anna}), \text{surname}(\text{gialli}), \text{status}(\text{professor}), \text{phone}(4444), \text{teaching}(\text{course}(\text{algebra})))\}.$$

Now, by applying Definition 12, we detect that Web page (1) is incomplete w.r.t. $\mu_1(\mathbf{e}_1)$, which is therefore an incompleteness symptom for (1). In fact, Web page (1) lacks teaching information.

Consider now the requirement $\mu_2(\mathbf{e}_2) \equiv (\text{member}(\text{name}(\text{anna}), \text{surname}(\text{gialli})))$. The associate likely missed information set is

$$\text{LMIS}_{\mu_2(\mathbf{e}_2)} = \{\text{members}(\text{member}(\text{name}(\text{mario}), \text{surname}(\text{rossi}), \text{status}(\text{professor})), \text{member}(\text{name}(\text{franca}), \text{surname}(\text{bianchi}), \text{status}(\text{technician})), \text{member}(\text{name}(\text{giulio}), \text{surname}(\text{verdi}), \text{status}(\text{student})))\}.$$

Note that the Web page in $LMIS_{\mu_2(e_2)}$ is incomplete w.r.t. the requirement $\mu_2(e_2)$, which models the fact that **anna gialli** must be a member of the group. Finally, the remaining requirements do not give rise to detect further errors.

It is worth pointing out that our verification framework is able to detect both the erroneous Web pages and the cause of the detected errors (i.e., the so-called 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 incompleteness symptoms, since the latter provides the missing information which is needed to complete the erroneous Web pages.

7 Implementation

The basic methodology presented so far has been implemented in the preliminary prototype system VERDI (VERification and Rewriting for Debugging Internet sites), which is written in DrScheme v205 [12] and is publicly available together with a set of tests at <http://www.dimi.uniud.it/~demis/#software>.

The implementation consists of about 80 function definitions (approximately 1000 lines of source code). VERDI includes a parser for semistructured expressions and Web specifications, and several modules implementing the user interface, the partial rewriting mechanism and the verification technique. The system allows the user to load a Web site consisting of a finite set of semistructured expressions 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. The user interface is guided by textual menus, which are (hopefully) self-explaining.

We tested the system on several Web site examples which can be found at the URL address mentioned above. In each considered test case, we were able to detect the errors (i.e. missing and incomplete Web pages) efficiently. For instance, it took less than one second the verification of the Web site of Example 13 w.r.t the Web specification of the Example 8, producing error messages when necessary.

8 Conclusions

Conceiving and maintaining Web sites is a difficult task. In this paper, we provide a rewriting-based, formal specification language which can be used to impose properties both on the structure (syntactic properties) and on the contents (semantic properties) of Web sites. The computation mechanism underlying this language is based on a novel rewriting-like technique, called *partial rewriting*, in which the traditional pattern matching mechanism is replaced with tree *simulation* [11]. In our methodology, Web sites are automatically checked w.r.t. a given Web specification in order to detect incomplete and missing Web pages. Moreover, by analyzing the requirements not fulfilled by the Web site, we are also able to find out the missing information needed to repair the Web site. Our methodology exploits some marking information on terms which represent the

Web pages to better locate the errors, which is provided by the user in advance. We have also discussed some implementation details of the preliminary system VERDI, a prototype implementation of the verification framework that we propose. Thus, we use the rewriting-based machinery as a common formalism for the specifications, for specifying our verification technique and for implementing the verification tool.

Finally, 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. We also plan to extend the specification language in order to support the detection of regular expressions. This is useful to guarantee that proprietary or “forbidden” data are not displayed on the external version of the site (e.g. a number of credit card). On the practical side, we plan to develop a fully user-friendly system which can help Web administrators to design, check and maintain their Web sites.

References

1. S. Abiteboul, P. Buneman, and D. Suciu. *Data on the Web. From Relations to Semistructured Data and XML*. Morgan Kaufmann, 2000.
2. M. Bezem. *TeReSe, Term Rewriting Systems*, chapter Mathematical background (Appendix A). Cambridge University Press, 2003.
3. F. Bry and S. Schaffert. Towards a Declarative Query and Transformation Language for XML and Semistructured Data: Simulation Unification. In *Proc. of the Int'l Conference on Logic Programming (ICLP'02)*, volume 2401 of *LNCS*. Springer-Verlag, 2002.
4. F. Bry and S. Sebastian Schaffert. The XML Query Language Xcerpt: Design Principles, Examples, and Semantics. Technical report, 2002. Available at: <http://www.xcerpt.org>.
5. P. Buneman, S. B. Davidson, G. G. Hillebrand, and D. Suciu. A Query Language and Optimization Techniques for Unstructured Data. In *Proc. of ACM SIGMOD Int'l Conference on Management of Data (ICMD'96)*, 1996.
6. A. Cortesi, A. Dovier, E. Quintarelli, and L. Tanca. Operational and abstract semantics of a graphical query language. *Theoretical Computer Science*, 275:521–560, 2002.
7. N. Dershowitz and D. Plaisted. Rewriting. *Handbook of Automated Reasoning*, 1:535–610, 2001.
8. T. Despeyroux and B. Trousse. Semantic Verification of Web Sites Using Natural Semantics. In *Proc. of 6th Conference on Content-Based Multimedia Information Access (RIAO'00)*, 2000.
9. M. Fernandez, D. Florescu, A. Levy, and D. Suciu. Verifying Integrity Constraints on Web Site. In *Proc. of Sixteenth International Joint Conference on Artificial Intelligence (IJCAI'99)*, volume 2, pages 614–619. Morgan Kaufmann, 1999.
10. M. F. Fernandez and D. Suciu. Optimizing regular path expressions using graph schemas. In *Proc. of Int'l Conference on Data Engineering (ICDE'98)*, pages 14–23, 1998.
11. M. R. Henzinger, T. A. Henzinger, and P. W. Kopke. Computing simulations on finite and infinite graphs. In *IEEE Symposium on Foundations of Computer Science*, pages 453–462, 1995.

12. PLT. DrScheme web site. Available at: <http://www.drscheme.org>.
13. World Wide Web Consortium (W3C). HyperText Markup Language (HTML) 4.01, 1997. Available at: <http://www.w3.org>.
14. World Wide Web Consortium (W3C). Extensible Markup Language (XML) 1.0, second edition, 1999. Available at: <http://www.w3.org>.

A Proof of the technical results

This appendix contains the proofs of the technical results of the paper.

Lemma 1. *Let $\mathbf{s}_1 \equiv (V_1, E_1, r_1, label_1)$, $\mathbf{s}_2 \equiv (V_2, E_2, r_2, label_2)$ be two Web page templates in $\tau(\mathcal{Text} \cup \mathcal{Tag}, \mathcal{V})$. A minimal page simulation \mathbf{S} of \mathbf{s}_1 in \mathbf{s}_2 w.r.t. \sim is a total mapping $\mathbf{S} : V_1 \rightarrow V_2$.*

Proof. Let $\mathbf{s}_1 \equiv (V_1, E_1, r_1, label_1)$, $\mathbf{s}_2 \equiv (V_2, E_2, r_2, label_2) \in \tau(\mathcal{Text} \cup \mathcal{Tag}, \mathcal{V})$. From points 1 and 3 of Definition 3, and by using the fact that \mathbf{s}_1 has an underlying tree structure (in particular, it is a connected graph), we have that

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

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

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

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

Let us first give a technical lemma in order to demonstrate Proposition 1.

Lemma 2. *Let \mathbf{I} be a Web specification and \mathbf{W} be a Web site. If $\mu(\mathbf{s}) \in \text{lfp}(\mathbf{R}_{\mathbf{I}})$, then \mathbf{s} is an instance of some $\mathbf{r} \in \text{Rhs}_{\mathbf{I}}$.*

Proof. Let $\mu(\mathbf{s}) \in \text{lfp}(\mathbf{R}_{\mathbf{I}})$. Therefore, from Definition 8, $\mu(\mathbf{s})$ is derived from some Web page $\mathbf{p} \in \mathbf{W}$ by using partial rewriting. Hence, we have $\varepsilon(\mathbf{p}) \rightarrow \dots \rightarrow \mu(\mathbf{s})$. From point (2) of Definition 7, $\mathbf{s} \equiv \mathbf{r}\sigma$, for some substitution σ and $\mathbf{r} \in \text{Rhs}_{\mathbf{I}}$. Thus, \mathbf{s} is an instance of some $\mathbf{r} \in \text{Rhs}_{\mathbf{I}}$.

We define the *height* of a Web page template $\mathbf{s} = (V, E, r, label)$ as $\max\{\text{depth}(\mathbf{s}, v) \mid v \in V\}$.

Proposition 1. *Let \mathbf{I} be a bounded Web specification and \mathbf{W} be a Web site. Then, there exists $k \in \mathbb{N}$ such that $\text{lfp}(\mathbf{R}_{\mathbf{I}}) = \mathbf{R}_{\mathbf{I}} \uparrow^{\mathbf{W}} k$.*

Proof (Sketch). By contradiction, let us suppose that there does not exist $k \in \mathbb{N}$ such that $\text{lfp}(\mathbf{R}_{\mathbf{I}}) = \mathbf{R}_{\mathbf{I}} \uparrow^{\mathbf{W}} k$. That is, $\text{lfp}(\mathbf{R}_{\mathbf{I}}) = \mathbf{R}_{\mathbf{I}} \uparrow^{\mathbf{W}} \omega$ where ω is the first infinite ordinal and, consequently, $\text{lfp}(\mathbf{R}_{\mathbf{I}})$ is infinite. By Lemma 2, there are infinite distinct elements in $\text{lfp}(\mathbf{R}_{\mathbf{I}})$ which are instances of a right-hand side \mathbf{r} of a rule of \mathbf{I} ; since Web site \mathbf{W} is finite, those instances can be generated only by partial rewriting. Then, consider an infinite partial rewrite sequence \mathcal{D} from some $\varepsilon(\mathbf{p}) \in \mathbf{W}$. \mathcal{D} contains an infinite number of distinct marked instances of the right-hand side \mathbf{r} . This implies that the height of these instances of \mathbf{r} is not bounded. Hence, this amounts to saying that \mathbf{r} is simulated by a $\mathbf{l} \in \text{Lhs}_{\mathbf{I}}$ and there is at least one vertex v of \mathbf{r} labeled with a variable such that $\text{depth}(\mathbf{r}, v) > \text{depth}(\mathbf{l}, v')$ for every vertex v' of \mathbf{l} with a variable as label. Thus, \mathbf{I} is not bounded, which contradicts the hypothesis.

The following definition is auxiliary.

Definition 13. Let \mathbf{S} be a page simulation of a Web page template $\mathbf{s} \equiv (V, E, r, \text{label})$ in a Web page template \mathbf{t} w.r.t \sim , and V' be a set of vertices such that $V' \subseteq V$. Then, a restriction $\mathbf{S}_{|V'}$ of \mathbf{S} w.r.t V' is defined as $\mathbf{S}_{|V'} = \{(v, v') \in \mathbf{S} \mid v \in V'\}$.

Proposition 2. Let \mathbb{W} be a Web site, \mathbb{I} be a bounded Web specification and $\text{Req}_{\mathbb{I}, \mathbb{W}}$ be the set of requirements for \mathbb{W} w.r.t. \mathbb{I} . Let $\mu(\mathbf{e}) \in \text{Req}_{\mathbb{I}, \mathbb{W}}$. If $\text{LMIS}_{\mu(\mathbf{e})} = \emptyset$, then for all $\mathbf{p} \equiv (V, E, r, \text{label}) \in \mathbb{W}$, and $v \in V$, there is no minimal injective simulation of \mathbf{e} in $\mathbf{p}_{|v}$.

Proof. Let $\mathbf{p} \equiv (V, E, r, \text{label}) \in \mathbb{W}$. By contradiction, let us suppose that there exists a vertex $v \in V$ and a minimal, injective simulation \mathbf{S} of $\mathbf{e} = (V_{\mathbf{e}}, E_{\mathbf{e}}, r_{\mathbf{e}}, \text{label}_{\mathbf{e}})$ in $\mathbf{p}_{|v}$ w.r.t \sim . Now, let us consider the marked part of \mathbf{e} , that is, $\text{mark}(\mu, \mathbf{e}) = (V_{\mu(\mathbf{e})}, E_{\mu(\mathbf{e})}, r_{\mathbf{e}}, \text{label}_{\mathbf{e}})$. Then, trivially, $V_{\mu(\mathbf{e})} \subseteq V_{\mathbf{e}}$, $E_{\mu(\mathbf{e})} \subseteq E_{\mathbf{e}}$. Thus, by Definition 3, $\mathbf{S}_{|V_{\mu(\mathbf{e})}}$ is a minimal, injective simulation of $\text{mark}(\mu, \mathbf{e})$ in $\mathbf{p}_{|v}$ w.r.t \sim . This implies that $\text{LMIS}_{\mu(\mathbf{e})} \neq \emptyset$, which is a contradiction. Therefore, there exists no minimal, injective simulation \mathbf{S} of $\mathbf{e} = (V_{\mathbf{e}}, E_{\mathbf{e}}, r_{\mathbf{e}}, \mathbf{e})$ in $\mathbf{p}_{|v}$ w.r.t \sim . \square