

Automatic Verification of Web Sites Using Partial Rewriting

M. Alpuente¹ D. Ballis² M. Falaschi²

¹ DSIC, Universidad Politécnica de Valencia, Spain.

`alpuente@dsic.upv.es`

² DIMI, Università degli Studi di Udine, Italy.

`{demis,falaschi}@dimi.uniud.it`

Talk plan

- Motivations
- Web site denotation and Web site specifications
- Simulation and partial rewriting
- Verification Framework
- Conclusions

Motivation

- Web sites can have a very complex structure.
- Web sites are difficult to **develop** and to **maintain**.

Motivation

- Web sites can have a very complex structure.
- Web sites are difficult to **develop** and to **maintain**.
- Suitable methodologies are needed to verify Web sites w.r.t a given scheme, which is able to express semantic as well as syntactic properties.
- **Formal methods** can be fruitfully employed!

Motivation

- Web sites can have a very complex structure.
- Web sites are difficult to **develop** and to **maintain**.
- Suitable methodologies are needed to verify Web sites w.r.t a given scheme, which is able to express semantic as well as syntactic properties.
- **Formal methods** can be fruitfully employed!
- We can verify and improve correctness of the EU-India project outcomes (Web sites)
- A workshop organized in Valencia in February 2005.
- Collaboration with Hyderabad on developments of new tools and improvements of basic framework.

Formal Verification of Web sites

We provide a rule-based specification language for

- **specifying** integrity conditions for a given Web site.

Formal Verification of Web sites

We provide a rule-based specification language for

- **specifying** integrity conditions for a given Web site.

And a verification technique for

- automatically **checking** whether those conditions are fulfilled.
- helping to **repair** faulty Web sites.

Our framework is based on a rewriting-like technique (which is called **partial rewriting**), more suitable for dealing with semistructured data (e.g. XML and/or HTML documents).

Web site denotation

In our framework, a Web page is a ground term. Consequently, we represent a Web site as a finite collection of **ground terms** of a suitable term algebra $\tau(\mathcal{T}ext \cup \mathcal{T}ag)$.

Example:

```
<members>
  <member>
    <name> Mario </name>
    <surname> Rossi </surname>
    <status> Professor </status>
  </member>
  <member>
    <name> Franca </name>
    <surname> Bianchi </surname>
    <status> Technician </status>
  </member>
</members>
```

```
members(
  member(
    name(Mario),
    surname(Rossi),
    status(Professor)
  ),
  member(
    name(Franca),
    surname(Bianchi),
    status(Technician)
  )
)
```


Specification Language

- Web Specifications are sets of rules which allow to specify conditions in order to
 - detect **forbidden** or **incorrect** information
 - detect **missing** or **incomplete** Web pages
- a Web specification is made up of
 - a set of **correctness rules** I_N
 - a set of **completeness rules** I_M
 - a set of **rewrite rules** (i.e. a Term Rewriting System) R

Correctness Rules

- They allow to detect an incorrect/forbidden pattern inside a Web page.
- A *correctness rule* has the following form

$$1 \rightarrow \text{error} \mid X_1 \text{ in } \text{rexp}_1 \dots X_n \text{ in } \text{rexp}_n$$

with $X_i \in \text{Var}(1)$ and rexp_i a regular expression over $\tau(\text{Text} \cup \text{Tag})$, $i = 1, \dots, n$.

- **Interpretation:** if 1 is recognized in some Web Page, then 1 is rewritten to `error` whenever the values bound to each X_i belong to the corresponding regular language rexp_i .

Correctness Rules - Examples

- A rule to forbid sexual contents from being published in the home pages.

$$\text{hpage}(X) \rightarrow \text{error} \mid X \text{ in } [:\text{TextTag}:] * \text{sex} [:\text{TextTag}:]*$$

- A rule to detect blinking text (aim: improving accessibility for people with disabilities)

$$\text{blink}(X) \rightarrow \text{error}$$

Completeness Rules

- They allow to discover missing information in a Web site.
- A completeness rule is a rule of the following form

$$l \rightarrow \mu(r) \langle q \rangle$$

where l, r are terms, μ is a marking function for marking some symbols of r by means of a $\#$.

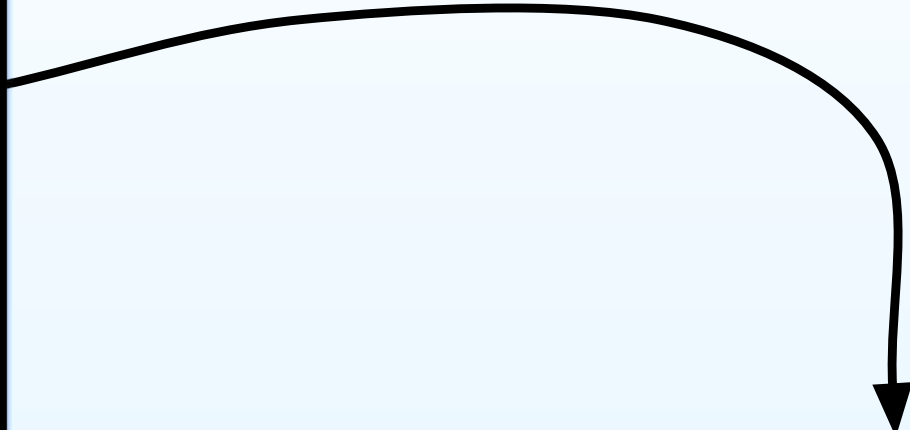
- $q \in \{A, E\}$, q can be either an **universal** or an **existential** quantifier.
- Marks are used to select the Web pages on which we want to check a given condition.

Marked part of a term

```
#hpage (  
  fullname(  
    "Mario Rossi"  
  )  
  #status(  
    #"Professor"  
  )  
  teaching(  
    course("Logic"),  
    course("Algebra")  
  )  
)
```

Marked part of a term

```
#hpage (  
  fullname (  
    "Mario Rossi"  
  )  
  #status (  
    #"Professor"  
  )  
  teaching (  
    course ("Logic"),  
    course ("Algebra")  
  )  
)
```



```
#hpage (  
  #status (  
    #"Professor"  
  )  
)
```

Completeness Rules - Interpretation

Given a Web site \mathbb{W}

- An **existential** completeness rule $\perp \multimap \mu(\mathbf{r})\langle\mathbf{E}\rangle$ is interpreted as follows:
 - if \perp is recognized in some Web page of \mathbb{W} , then (the irreducible form of) \mathbf{r} must be recognized in **some** Web page of \mathbb{W} which contain the marked part of \mathbf{r} .
- An **universal** completeness rule $\perp \multimap \mu(\mathbf{r})\langle\mathbf{A}\rangle$ is interpreted as follows:
 - if \perp is recognized in some Web page of \mathbb{W} , then (the irreducible form of) \mathbf{r} must be recognized in **every** Web page of \mathbb{W} which contain the marked part of \mathbf{r} .

Completeness Rules - Examples

- A rule to verify whether the home pages of professors contain some teaching information inside.

$$\text{hpage}(\text{status}(\text{professor})) \rightarrow \#\text{hpage}(\#\text{status}(\#\text{professor}), \text{teaching}) \langle \text{A} \rangle$$

- A rule to verify whether there exists a home page for each member of a group.

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

- A rule to verify whether each author of a publication is also a member of the group.

$$\text{pubs}(\text{pub}(\text{name}(X), \text{surname}(Y))) \rightarrow \#\text{member}(\text{name}(X), \text{surname}(Y)) \langle \text{E} \rangle$$

A Web specification - Example

Consider a Web site containing some information about a research group (e.g., member group affiliation, personal data, publications, ...).

- Correctness Rules I_N

$$\text{hpage}(X) \rightarrow \text{error} \mid X \text{ in } [:\text{TextTag}:] * \text{sex } [:\text{TextTag}:]*$$
$$\text{blink}(X) \rightarrow \text{error}$$

- Completeness Rules I_M

$$\text{hpage}(\text{status}(\text{professor})) \rightarrow \#\text{hpage}(\#\text{status}(\#\text{professor}), \text{teaching}) \langle A \rangle$$
$$\text{member}(\text{name}(X), \text{surname}(Y)) \rightarrow \#\text{hpage}(\text{fullname}(\text{append}(X, Y)), \text{status}) \langle E \rangle$$
$$\text{pubs}(\text{pub}(\text{name}(X), \text{surname}(Y))) \rightarrow \#\text{member}(\text{name}(X), \text{surname}(Y)) \langle E \rangle$$

- Rewrite Rules R = Definition of function `append`

Tree Simulation

Simulation allows us to recognize the structure and the labels of a Web page (template) into another. It provides a powerful pattern-matching mechanism:

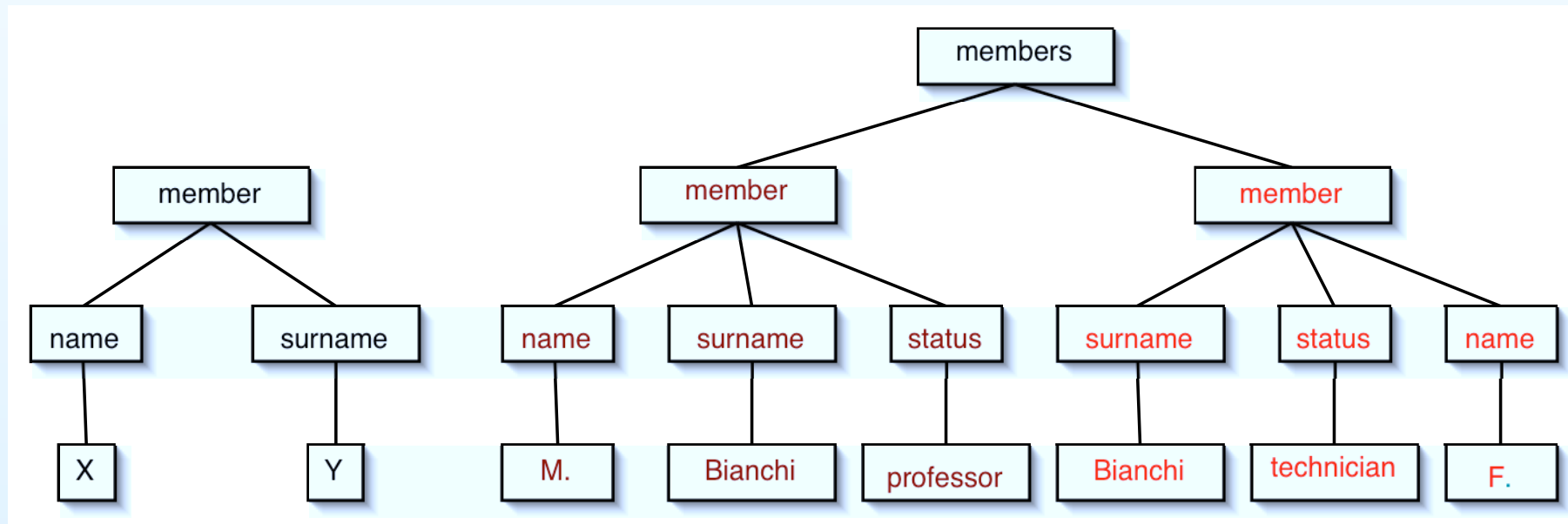
- suitable for dealing with HTML/XML data (partial matching, unordered trees)
- fast (efficient algorithms do exist)

Tree Simulation

Simulation allows us to recognize the structure and the labels of a Web page (template) into another. It provides a powerful pattern-matching mechanism:

- suitable for dealing with HTML/XML data (partial matching, unordered trees)
- fast (efficient algorithms do exist)

Example (minimal, injective simulation)

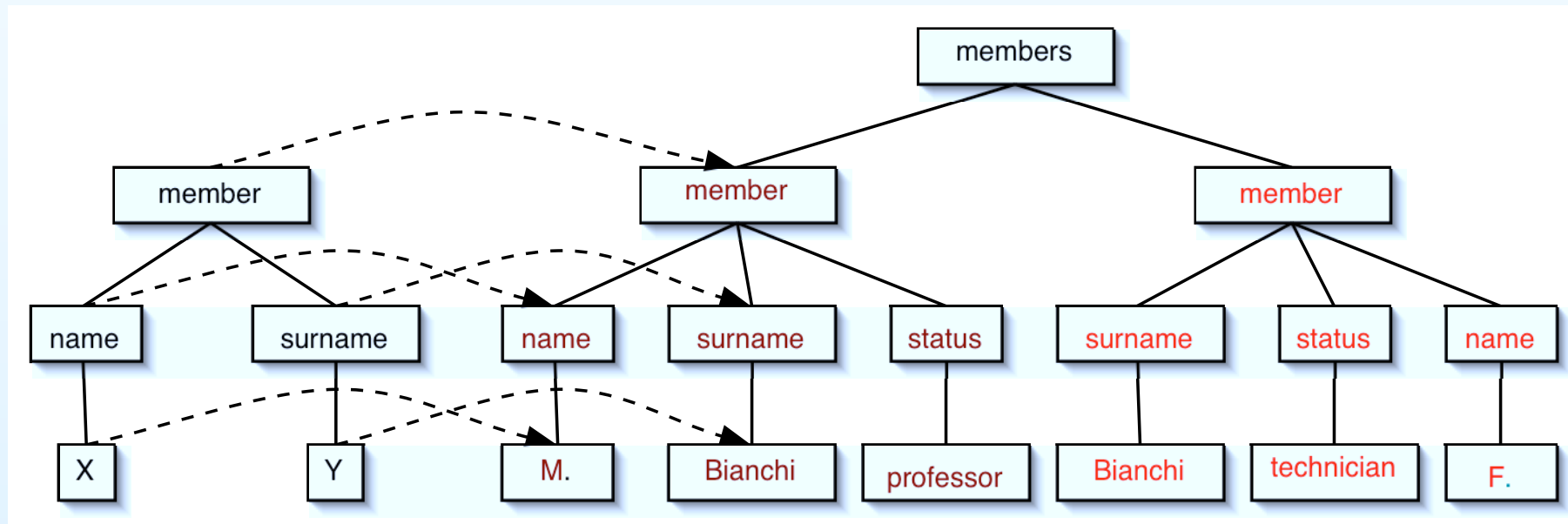


Tree Simulation

Simulation allows us to recognize the structure and the labels of a Web page (template) into another. It provides a powerful pattern-matching mechanism:

- suitable for dealing with HTML/XML data (partial matching, unordered trees)
- fast (efficient algorithms do exist)

Example (minimal, injective simulation)

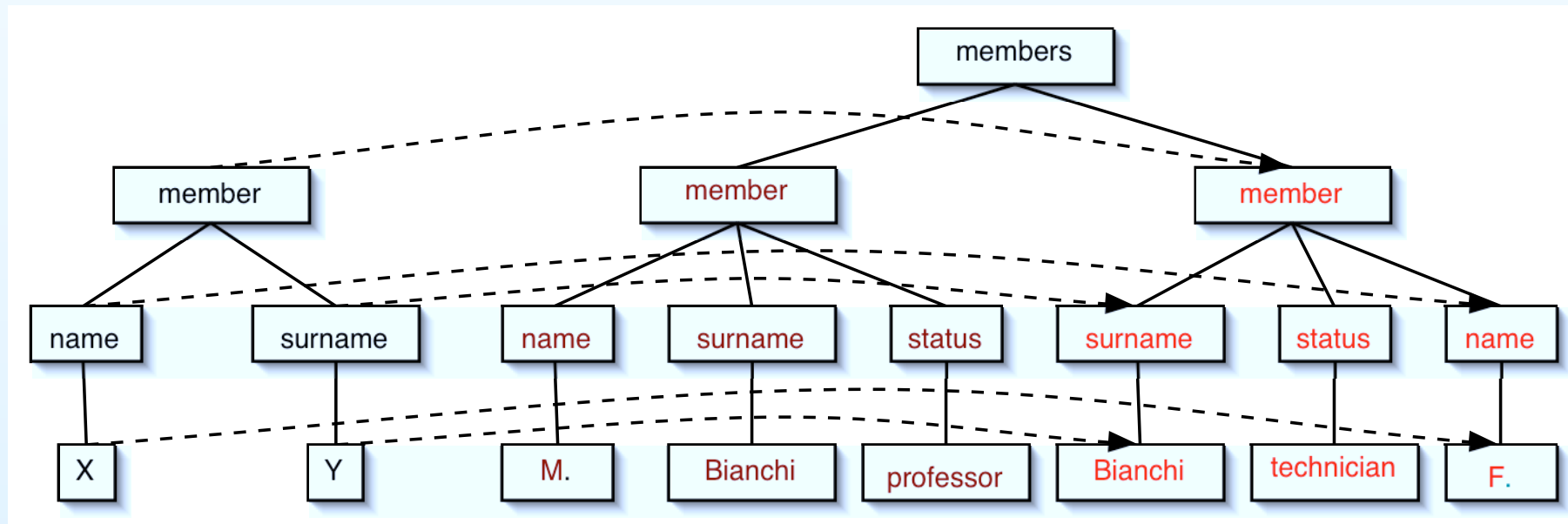


Tree Simulation

Simulation allows us to recognize the structure and the labels of a Web page (template) into another. It provides a powerful pattern-matching mechanism:

- suitable for dealing with HTML/XML data (partial matching, unordered trees)
- fast (efficient algorithms do exist)

Example (minimal, injective simulation)



Partial Rewriting \rightarrow

Partial Rewriting: a rewriting relation in which

- the traditional *pattern matching* mechanism is replaced by tree **simulation**
- the context of selected reducible expressions is disregarded
- we deal with marking information

$$\mu(\mathbf{t}) \rightarrow_r^\sigma \mu'(\mathbf{t}')$$

Partial rewriting steps

```
members(member(name(Mario), surname(Rossi),
               status(Professor)),
         member(name(Franca), surname(Bianchi),
               status(technician))
        )
```

is partially rewritten to

```
hpage(fullname(append(Mario,Rossi), status)
→R hpage(fullname(MarioRossi), status)
```

and

```
hpage(fullname(append(Franca,Bianchi), status)
→R hpage(fullname(FrancaBianchi), status)
```

by rule $\text{member}(\text{name}(X), \text{surname}(Y)) \rightarrow \underline{\text{hpage}}(\text{fullname}(\text{append}(X, Y)), \text{status})$

Verification of a Web site

We provide a verification methodology which is able to

1. detect erroneous patterns of a Web sites \bar{w} w.r.t a Web specification \mathcal{I} ;
2. detect missing/incomplete Web pages of a Web sites \bar{w} w.r.t a Web specification \mathcal{I} ;
3. collect information which is needed to repair \bar{w} .

Erroneous pattern detection

Methodology

1. Check whether some Web page of W can be **partially** rewritten to error by using some correctness rule r_{corr} .
2. Check whether the instantiated condition of r_{corr} holds.

If Points (1) and (2) are satisfied, an incorrectness error is raised!

Missing/incomplete Web pages

The main idea is the following.

1. Generate a set of requirements $Req_{I,W}$ by applying I to W .
2. Check whether requirements in $Req_{I,W}$ are fulfilled in W .

Terms, which are derived via partial rewriting from the original Web site by means of the completeness rules, can be considered as **requirements** to be fulfilled by W .

We can compute the set of requirements via a fixpoint computation.

$$Req_{I,W} = lfp(R_I)$$

where R_I is a continuous operator similar to the immediate consequence operator which is used in Logic Programming.

Missing/incomplete Web pages (cont'd)

A requirement **is not fulfilled** if it cannot be simulated in a suitable subset of the Web site which has been generated by exploiting marking information.

Existential and universal checks!

Repairing Web sites

During the verification process, we can collect useful information in order to semi-automatically fix the correctness/completeness errors which are discovered. That is,

- detected incorrect patterns allow us to locate the information we need to change in order to get a correct Web site.
- unsatisfied requirements provide us the missing information which is necessary to complete the Web site.

Conclusions

- We presented a rewriting-based language for the specification and the verification of integrity conditions on Web sites.
 - partial rewriting
 - erroneous patterns and missing/incomplete Web pages
- We have also implemented the prototype VERDI (VERification and Rewriting for Debugging Internet sites), which is publicly available at

`http://www.dimi.uniud.it/~demis/#software`

- We tested the system on real Web sites, e.g.

`http://www.dimi.uniud.it/clg`

END