

# Computational Aspects of Resilient Data Extraction from Semistructured Sources\*

(Extended Abstract)

Hasan Davulcu    Guizhen Yang    Michael Kifer    I.V. Ramakrishnan

Department of Computer Science  
SUNY at Stony Brook  
Stony Brook, NY 11794-4400, USA

{davulcu,guizyang,kifer,ram}@cs.sunysb.edu

## ABSTRACT

Automatic data extraction from semistructured sources such as HTML pages is rapidly growing into a problem of significant importance, spurred by the growing popularity of the so called "shopbots" that enable end users to compare prices of goods and other services at various web sites without having to manually browse and fill out forms at each one of these sites.

The main problem one has to contend with when designing data extraction techniques is that the contents of a web page changes frequently, either because its data is generated dynamically, in response to filling out a form, or because of changes to its presentation format. This makes the problem of data extraction particularly challenging, since a desirable requirement of any data extraction technique is that it be "resilient", i.e., using it we should always be able to locate the object of interest in a page (such as a form or an element in a table generated by a form fill-out) in spite of changes to the page's content and layout.

In this paper we propose a formal computation model for developing resilient data extraction techniques from semistructured sources. Specifically we formalize the problem of data extraction as one of generating unambiguous extraction expressions, which are regular expressions with some additional structure. The problem of resilience is then formalized as one of generating a maximal extraction expression of this kind. We present characterization theorems for maximal extraction expressions, complexity results for testing them, and algorithms for synthesizing them.

---

\*Work supported in part by the NSF grants CCR-9711386 and EIA-9705998.

## 1. INTRODUCTION

The World Wide Web is becoming the dominant medium for information delivery and electronic commerce. The number of users who routinely use the web to buy goods and services continues to increase at a rapid pace. In response, software robots (called "shopbots") that allow consumers to quickly find out the best prices for comparable goods and services are beginning to emerge. Information about prices and other attributes of products are typically obtained by filling out forms at a vendor's site. Software robots retrieve such information by automatically navigating to relevant sites, locating the correct forms, filling them out and extracting the data of interest from web pages returned as the result. (Junglee and Jango [17] are two examples of such shopbots.) Currently almost all web information is stored as semistructured data, mostly as HTML pages, and shopbots need to retrieve data from such sources. Hence automatic data extraction from semistructured data sources is a problem of significant importance especially in the context of web-based electronic commerce.

This problem has attracted a lot of research attention recently. The techniques proposed so far, are by and large centered around creating a *wrapper* [3, 6, 9, 12, 13, 15, 18, 14, 7, 21, 22], that parses an HTML source and maps it into a set of structured or semistructured database objects that can be readily queried and manipulated by applications.

The central difficulty in designing data extraction techniques is the volatile nature of HTML pages, in the sense that they change very frequently. Changes occur either because they have to accommodate new services and content offerings or because they are *dynamically* generated in response to form-based queries. Such variations can give rise to "brittleness" in data extractors, i.e. they may no longer be able to locate the objects of interest in the page (such as a form or a table element). Thus developing data extraction techniques that are resilient to changes in the data source structure is both desirable and important. To the best of our knowledge this problem has not yet been fully explored.

In this paper, we make the first step towards developing a formal framework for creating resilient data extraction wrappers for semistructured sources. As an example, Web pages

can be represented as sequences of tokens (HTML tags and strings), and the extraction problem is usually reduced to the problem of parsing using regular grammars [9, 12, 13, 21, 18, 22, 14, 15], context-free grammars [6, 3] or specialized languages [7].

We propose the notion of *extraction expressions*, which are tag-marked regular expressions, as a formalization of the informal concept of the “target object that can be identified by its local or global context.”

The initial extraction expression is usually derived from a set of examples (i.e., of HTML pages and target objects) using one of the known heuristic approaches [13, 21, 22, 9], or one of the learning algorithms [18, 4, 5], or even manually [14, 15, 12].

We first introduce the notion of “unambiguity” as a consistency requirement for extraction expressions. Unambiguous extraction expressions must identify the target objects uniquely within any page. We show that ambiguity of any given extraction expression is decidable in polynomial time.

If an extraction expression can correctly identify the target object in many variations of the same document, we shall call such an expression *resilient* to changes. Unambiguous extraction expressions can be ordered according to the languages which they parse. The larger the language — the more resilient the expression. We thus reduce the problem of resilience against variations to the task of synthesizing *maximal* unambiguous extraction expressions.

We show that some unambiguous expressions can be maximized in several different — even infinite number (!) of — ways and the problem of deciding maximality is PSPACE-complete. We do not know whether every unambiguous extraction expression can be maximized (although we conjecture that this is possible). However, we propose algorithms for maximizing a large, non-trivial class of extraction expressions that arise frequently in practical situations. Our preliminary experiments with the maximization algorithms proposed here indicate that they are sufficient to provide resilient extraction capabilities for a web-based information harvesting system that we have developed recently.

The remainder of the paper is organized as follows. Section 2 surveys the related work. In Section 3, we motivate the problem of resilient extraction and our approach using a concrete and realistic example. Section 4 introduces the main definitions, including the notions of ambiguity and resilience. In Section 5, we resolve the complexity of various decision problems related to resilient extraction. In Section 6, we present our maximization algorithms and prove their correctness. Section 7 illustrates how the techniques developed in this paper apply to the example of Section 3. Section 8 concludes our paper. Due to space limitation, we cannot present the proofs here. However, all proofs can be found in a technical report at <http://www.cs.sunysb.edu/~guizyang/papers/extraction.ps>

## 2. RELATED WORK

Semi-structured data [10, 1] has recently emerged as an important area of study. Querying semistructured data through

a query language [11, 2] is one way to extract data. This technique is applicable when the structure of the data and the location of the desired object is known to a large degree and is not subject to drastic changes.

Wrapper based extraction techniques from semistructured data include [3, 6, 9, 13, 15, 18, 22, 14, 7, 21]. The works [14, 15] provide a powerful framework for manual extraction of complex data objects based on the application of sequences of patterns. An expressive specialized language for data extraction is proposed in [7]. The language is based on searches using a subset of regular expressions and other page restructuring instructions. The works in [3, 9, 13, 20, 22] propose various graphical tools and heuristics for semi-automatic derivation of regular extraction expressions. Such techniques could be used at the initial learning stage of our framework (see Sections 3 and 7). However, none of these works addresses the problems inherent in resilient data extraction.

The works [18, 21] describe fully automated wrapper generation techniques that use machine learning induction algorithms (similar to [8, 4, 5]). Such techniques are useful in our framework since they could supply us with initial extraction expressions that could then be generalized using our techniques. A limitation of these approaches, however, is that the extracted data must be representable as a set of tuples. An ontology-based approach to extracting data is presented in [12]. This technique requires that suitable ontologies of context keywords, relationships, and regular expressions for lexicons must be constructed in advance (and manually). In [6], a heuristic algorithm for automatic generation of context-free grammar for data extraction purposes is presented. However this work does not address the issue robustness of the wrappers in the presence of variations in source documents.

Overall, with the exception of [6], all data extraction techniques we are aware of rely on regular expressions. Hence, our results enhance these techniques by providing both the objective criteria for sorting out the good expressions from the bad ones, and techniques for making the good extraction expressions more robust.

## 3. MOTIVATING EXAMPLE

One of the main problems with data extraction from Web documents is that the structure of the data might change due to the dynamic nature of the document or because the document was redesigned manually and, thus, the location of the object of interest might be hard to pinpoint. The most typical changes are insertion or deletion of HTML elements before or after the object of interest and embedding of the object inside some other HTML element.

To illustrate, consider the documents in Figures 1 and 2, which represent the main elements of a typical catalog shopping site. The user can fill out a form to search for a product by keywords or part number and out come the results. The top document consists of some header information and a form. The bottom document is similar, but the form is now embedded in a table and a new table row, a link to customer service, is added. If we are building a web robot that goes around and compares prices, the most likely object of inter-

```

<P>
<H1><IMG SRC="supplier.gif">
  Virtual Supplier, Inc.
</H1>
<P>
<FORM METHOD="POST" ACTION="search.cgi">
<INPUT TYPE="image" ALIGN=left
  SRC="search.gif">
<INPUT TYPE="text" SIZE="15"
  NAME="VALUE"> <BR>
<P>
<INPUT TYPE="radio" NAME="ATTR"
  VALUE="1" checked>
Keywords<BR>
<INPUT TYPE="radio" NAME="ATTR"
  VALUE="2">
Manufacturer Part#
</FORM>

```

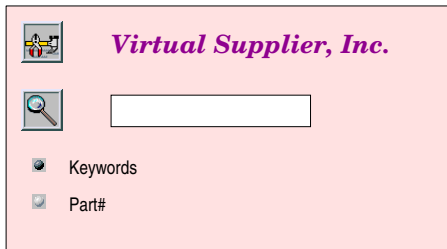


Figure 1: Original Page

est would be the form. Suppose now that we train our robot to find the needed form using the top document. We want to make sure that the robot does not fail even if the site administrator reorganizes the document as shown in Figure 2, if more rows are added to the second document before or after the form, and (hopefully) if other forms or tables are added to that document.

To begin, we can try to represent our documents as strings of objects. For instance, the top document can be represented as

P H1 /H1 P FORM /FORM

In this representation, we assume that the contents of the objects H1 and FORM is of no interest. If the insides of, say, the form are of interest (usually they are), we could expand the representation as follows:

P H1 /H1 P FORM INPUT INPUT P INPUT INPUT /FORM

Likewise, the second document in the figure can be represented as:

```

TABLE TR TD /TD TD /TD /TR TR TD /TD TD /TD /TR
FORM TR TD INPUT /TD TD INPUT /TD TD INPUT /TD /TR
/FORM /TABLE

```

The exact details of this representations is of no significance here. It is easy to enrich this model to take the tag attributes into account, and so on. The point is that documents can be represented as sequences of tags plus some additional

```

<TABLE>
<TR> <TD><IMG SRC="supplier.gif"></TD>
  <TD><H1>Virtual Supplier, Inc.
  </H1>
  </TD>
</TR>
<TR> <TD><IMG SRC="cust.gif"></TD>
  <TD><A HREF="cust.html">
    Customer Service
  </A><BR></TD>
</TR>
<FORM METHOD="POST" ACTION="search.cgi">
<TR> <TD><INPUT TYPE="image"
  SRC="search.gif">
  </TD>
  <TD><INPUT TYPE="text"
    SIZE="15" NAME="VALUE">
  </TD>
  <TD> <INPUT TYPE="radio" NAME="ATTR"
    VALUE="1" checked>
    Keywords<BR>
    <INPUT TYPE="radio" NAME="ATTR"
    VALUE="2">
    Manufacturer Part#
  </TD>
</TR>
</FORM>
</TABLE>

```

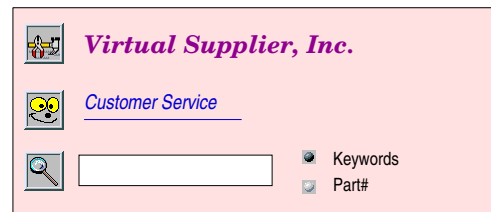


Figure 2: Rearranged Page

symbols (*e.g.*, to incorporate tag attributes). Each of the above strings is an abstract representation of a semistructured Web document, but alone they do not yet specify the objects of interest. To tell the robot which object we want, we can enclose it in angle brackets. To indicate that we are interested in the second INPUT-element of the form, we can write the following:

P H1 /H1 P FORM INPUT <INPUT> P INPUT INPUT /FORM

and

```

TABLE TR TD /TD TD /TD /TR TR TD /TD TD /TD /TR
FORM TR TD INPUT /TD TD <INPUT> /TD TD INPUT /TD /TR
/FORM /TABLE

```

Since the above expressions are obviously different and each matches only one of our candidate documents, they are still not very useful to our robot. However, we can try to generalize them, say, using the syntax of regular expressions:

```

(Tags - FORM)* FORM (Tags - INPUT)* INPUT
(Tags - INPUT)* <INPUT> Tags*

```

In the above expression,  $\mathcal{T}ags$  is a regular expression that matches any HTML tag,  $(\mathcal{T}ags - FORM)$  matches any tag ex-

cept FORM), and “\*” says that the corresponding expression can be repeated zero or more times.

Since this regular expression matches both documents and identifies the object of interest correctly, we could give it to our robot and it will extract the information properly even if the document changes as shown in Figure 2.

Regular expressions with a marked occurrence of a symbol will be called *extraction expressions* in this paper. If an extraction expression can correctly identify the desired object in many variations of the same document, we shall say that such an expression is *resilient* to changes.

The above discussion suggests the following strategy for finding resilient extraction expressions. In the first stage, a small number of sample variants of the desired document can be obtained by filling out the same form in different ways. If the document is static but might change as a result of a redesign, heuristics can be used to automatically create several perturbations of the document. Next, the variations of the document in question are represented as strings and the object of interest is marked in each document (this object must be of the same kind in each case, since we assumed that they are perturbations of the same document), which leaves us with a small number of very rigid extraction expressions. Finally, these expressions are generalized into a single extraction expression that matches all the instances of our document.

Once the problem of data extraction is reduced to the problem of generalizing marked regular expressions, many interesting questions arise. The first of these is: how far should (and can) we generalize before the robot gets confused? Indeed, the original extraction expressions were not resilient because they matched only one document. However, when they *did* match, the object of interest was identified correctly and uniquely. Can it happen that a generalized extraction expression might match two or more objects and thus confuse our robot (because it will not know which is the true object of interest)? We do not need to think hard to find such a generalization:  $\text{Tags* } \langle \text{INPUT} \rangle \text{ Tags*}$ .

This leads to the question of *ambiguity* of extraction expressions. For instance, the expression  $((\text{qp}) * | (\text{Tags} - \text{p}) * ) \langle \text{p} \rangle \text{p} *$ , where “|” denotes an alternative (union of regular languages) and  $(\text{Tags} - \text{p})$  matches anything but  $\text{p}$ , is *unambiguous* in the sense that whenever it matches a string, the marked occurrence of the symbol  $\text{p}$  falls into a unique place on the string (even though the prefix  $((\text{qp}) * | (\text{Tags} - \text{p}) * )$  can match the prefix of a string in more than one way). In contrast,  $(\text{qp}) * \text{p} * \langle \text{p} \rangle \text{p} *$  is an *ambiguous* expression, because there are many ways to match the marked occurrence of  $\text{p}$  on *some* strings. For instance, consider  $\text{qppp}$ . The prefix  $(\text{qp}) * \text{p} *$  can match  $\text{qp}$  and  $\text{qpp}$ , so the marked occurrence of  $\text{p}$  can match the second *or* the third occurrence of  $\text{p}$ .

We believe that (un)ambiguity is an important property of extraction expressions that places limits on how much we can generalize from a set of examples. We shall see that ambiguity can be decided in polynomial time.

While unambiguity is important, we should not lose sight

of the fact that our goal is to find generalizations that are as resilient as possible. It turns out that extraction expressions of the form  $E_1 \langle \text{p} \rangle E_2$  can be ordered according to how vast the regular languages  $L(E_1)$  and  $L(E_2)$  are. The bigger they are — the more resilient the extraction expression is. Ideally, we would like to find generalizations that are *maximal* with respect to the above order among the unambiguous expressions that match the object of interest on our sample pages. We show that maximality of any given unambiguous regular expressions is decidable, albeit it is a PSPACE-complete problem. Surprisingly, the question of whether every unambiguous extraction expression has a maximal generalization seems to be a hard problem, and it remains open at this point. However, we shall see that maximal generalizations do exist for large classes of extraction expressions.

## 4. EXTRACTION EXPRESSIONS, UNAMBIGUITY, AND MAXIMALITY

The previous section provided the intuition behind the notions of extraction expressions, ambiguity, and resilience. We are now going to define these concepts formally and prove several of their important properties.

*Definition 4.1. (Extraction Expression)* Given a finite alphabet  $\Sigma$ ,  $p \in \Sigma$ , and regular expressions  $E_1$  and  $E_2$  over  $\Sigma$ ,  $E_1 \langle p \rangle E_2$  is called an extraction expression over  $\Sigma$ .

In other words, an extraction expression is simply an ordinary regular expression of a special form  $E_1 \langle p \rangle E_2$ , with one marked occurrence,  $\langle p \rangle$ , of an alphabet symbol.  $\square$

The language parsed by  $E_1 \langle p \rangle E_2$  is defined as  $L(E_1 \langle p \rangle E_2) \stackrel{\text{def}}{=} \{\rho \mid \rho \in L(E_1 \cdot p \cdot E_2)\}$ . Given a string  $\rho \in \Sigma^*$ , we say that  $E_1 \langle p \rangle E_2$  *parses*  $\rho$  if  $\rho \in L(E_1 \langle p \rangle E_2)$ . We also say that  $E_1 \langle p \rangle E_2$  *extracts*  $p$  from  $\rho$  if there exist  $\alpha, \beta \in \Sigma^*$ , such that  $\rho = \alpha \cdot p \cdot \beta$  and  $\alpha \in L(E_1), \beta \in L(E_2)$ .

Suppose  $\rho$  is a string. We can use  $E_1 \langle p \rangle E_2$  to extract information from  $\rho$  as follows: First we can try to split  $\rho$  into a prefix  $\alpha$ , followed by  $p$ , followed by a suffix  $\beta$ . If  $\alpha$  is recognized by  $E_1$  and  $\beta$  is recognized by  $E_2$ , then  $p$  is the extracted object. We try such splits until we either succeed on some split or fail on all candidates.

As illustrated in Section 3, the expression  $E_1 \langle p \rangle E_2$  might be *ambiguous*, i.e., there might be a string,  $\rho$ , where two different splits can succeed. For instance,  $\text{p} * \langle \text{p} \rangle \text{q}$  parses  $\text{pppq}$ , but any one of three  $\text{p}$ 's in  $\text{pppq}$  can be returned as the extracted object.

*Definition 4.2. (Unambiguous Extraction Expression)* Expression  $E_1 \langle p \rangle E_2$  is unambiguous iff for all  $\alpha, \alpha' \in L(E_1)$  and  $\beta, \beta' \in L(E_2)$ , if  $\alpha \cdot p \cdot \beta = \alpha' \cdot p \cdot \beta'$  then  $\alpha = \alpha'$  and  $\beta = \beta'$ .  $\square$

*Example 4.3. (Ambiguous and Unambiguous Extraction Expressions)*  $(\text{pq}) * \langle \text{p} \rangle \Sigma^*$  and  $(\text{p|pp}) \langle \text{p} \rangle (\text{p|pp})$  are ambiguous extraction expressions, whereas  $(\text{qp}) * \langle \text{p} \rangle \Sigma^*$  and  $(\text{p|pp}) \langle \text{p} \rangle (\text{p|ppp})$  are unambiguous extraction expressions. For

instance,  $pqpq$  can be parsed by  $(pq)^* \langle p \rangle \Sigma^*$  as  $\epsilon \cdot p \cdot qpq$  and as  $pq \cdot p \cdot q$ . Likewise,  $(p|pp)\langle p \rangle (p|pp)$  can parse  $pppp$  in two different ways. On the other hand, it can be proved that the last two expressions always parse their matching strings uniquely.  $\square$

Since, as explained in Section 3, we are mostly interested in unambiguous extraction expressions, from now on the term “extraction expression” will refer to unambiguous expressions, unless explicitly specified otherwise.

**Definition 4.4. (Partial Order among Extraction Expressions)**  $F_1 \langle p \rangle F_2 \preceq E_1 \langle p \rangle E_2$  iff  $L(F_1) \subseteq L(E_1)$  and  $L(F_2) \subseteq L(E_2)$ . We shall also say that  $E_1 \langle p \rangle E_2$  generalizes  $F_1 \langle p \rangle F_2$ .  $\square$

Informally, extraction expressions that are “larger” with respect to  $\preceq$  are more resilient because they can uniquely parse larger sets of strings. Moreover if  $F_1 \langle p \rangle F_2 \preceq E_1 \langle p \rangle E_2$  then the two expressions parse the strings in  $L(F_1 \langle p \rangle F_2)$  the same way. It is easy to see that  $\preceq$  is reflexive, antisymmetric and transitive. Therefore,  $\preceq$  is a partial order on the set of all unambiguous extraction expressions.

Note that  $F_1 \langle p \rangle F_2 \preceq E_1 \langle p \rangle E_2$  implies  $L(F_1 \langle p \rangle F_2) \subseteq L(E_1 \langle p \rangle E_2)$ , but not the other way around! Indeed, the two expressions  $p \langle p \rangle ppp$  and  $pp \langle p \rangle pp$  parse exactly the same language, but they extract different objects from that language:  $p \langle p \rangle ppp$  extracts the second occurrence of  $p$ , while  $pp \langle p \rangle pp$  extracts the third.

**Definition 4.5. (Maximal Extraction Expression)** An unambiguous extraction expression  $\mathcal{E}$  over a finite alphabet  $\Sigma$  is maximal iff for any unambiguous extraction expression  $\mathcal{E}'$  over  $\Sigma$ , if  $\mathcal{E} \preceq \mathcal{E}'$  then  $L(\mathcal{E}) = L(\mathcal{E}')$ .  $\square$

For the following examples, we use the expression  $E_1 - E_2$  to represent the regular expression that recognizes the regular set  $L(E_1) - L(E_2)$ .

**Example 4.6. (Maximal Extraction Expressions)** Although it might not be immediately obvious, both  $(\Sigma - p)^* \langle p \rangle \Sigma^*$  and  $(qp)^* \cdot ((\Sigma - p)^* - q) \langle p \rangle \Sigma^*$  are maximal extraction expressions.  $\square$

**Example 4.7. (Unambiguity and Maximality)** Given a non-maximal unambiguous extraction expression  $\mathcal{E}$  over  $\Sigma$ , if there exists a maximal extraction expression  $\mathcal{E}'$  over  $\Sigma$  such that  $\mathcal{E} \preceq \mathcal{E}'$ , we say that extraction expression  $\mathcal{E}$  can be maximized to  $\mathcal{E}'$ .

It is not known whether every unambiguous extraction expression  $\mathcal{E}$  can be maximized. However, even when maximization is known to exist then it might not be unique. For example,  $qp \langle p \rangle \Sigma^*$  can be maximized to

$$(\Sigma - p)^* \cdot p \cdot (\Sigma - p)^* \langle p \rangle \Sigma^*$$

and

$$((qp(\Sigma - p)^*) | ((\Sigma - p)^* - q)) \langle p \rangle \Sigma^*$$

(The latter expression is obtained using Algorithm 6.2 in Section 6, when  $qp \langle p \rangle \Sigma^*$  is used as input.) In fact, it can be shown that the above expression has an *infinite* number of maximal expressions that generalize it.  $\square$

## 5. COMPLEXITY OF THE AMBIGUITY AND MAXIMALITY PROBLEMS

**Definition 5.1. (Prefix and Suffix Factoring)** Given a pair of regular expressions  $E_1$  and  $E_2$  over a finite alphabet  $\Sigma$ , the prefix factorization of  $E_1$  by  $E_2$  is defined as  $E_2 \setminus E_1 \stackrel{\text{def}}{=} \{\alpha | \exists \beta \in L(E_2), \beta \cdot \alpha \in L(E_1)\}$ . The suffix factorization of  $E_1$  by  $E_2$  is  $E_1 / E_2 \stackrel{\text{def}}{=} \{\alpha | \exists \beta \in L(E_2), \alpha \cdot \beta \in L(E_1)\}$ .  $\square$

It is known that  $E_2 \setminus E_1$  and  $E_1 / E_2$  are regular languages, if both  $E_1$  and  $E_2$  are regular expressions [19]. Thus factors can be represented as regular expressions. Since every regular language corresponds to a regular expression and vice versa, we shall use  $E$  and  $L(E)$  interchangeably to denote the regular language recognized by  $E$ .

**LEMMA 5.2.** Given regular expressions  $E_1$  and  $E_2$  over a finite alphabet  $\Sigma$ , the regular expressions corresponding to  $E_2 \setminus E_1$  or  $E_1 / E_2$  can be computed in polynomial time in the size of  $E_1$  and  $E_2$ .

**LEMMA 5.3.**  $E_1 \langle p \rangle E_2$  is ambiguous iff there exist  $\alpha, \beta, \gamma \in \Sigma^*$  such that  $\alpha \cdot p \cdot \gamma \cdot p \cdot \beta \in L(E_1 \langle p \rangle E_2)$ ,  $\alpha, \alpha \cdot p \cdot \gamma \in L(E_1)$  and  $\beta, \gamma \cdot p \cdot \beta \in L(E_2)$ .

**PROPOSITION 5.4. (Necessary and Sufficient Condition for Unambiguity)** An extraction expression  $E_1 \langle p \rangle E_2$  over a finite alphabet  $\Sigma$  is unambiguous iff

$$(E_1 \cdot p) \setminus E_1 \cap E_2 / (p \cdot E_2) = \phi$$

**THEOREM 5.5. (Complexity of Testing Ambiguity)** Let  $E_1 \langle p \rangle E_2$  be an extraction expression over a finite alphabet  $\Sigma$ . Then testing whether  $E_1 \langle p \rangle E_2$  is ambiguous can be done in time quadratic in the size of  $E_1 \langle p \rangle E_2$ .

**PROPOSITION 5.6. (Necessary and Sufficient Condition for Maximality)** An unambiguous extraction expression  $E_1 \langle p \rangle E_2$  over a finite alphabet  $\Sigma$  is maximal iff

$$(E_1 \cdot p \cdot E_2) / (p \cdot E_2) = \Sigma^*$$

and

$$(E_1 \cdot p) \setminus (E_1 \cdot p \cdot E_2) = \Sigma^*$$

**LEMMA 5.7.** Given a regular expression  $E$  over a finite alphabet  $\Sigma$ , the problem of testing whether  $L(E) = \Sigma^*$  is PSPACE-complete.

**Proof:** See [16]  $\square$

LEMMA 5.8. *For any regular expression  $E$  over a finite alphabet  $\Sigma$ , the extraction expression  $(\Sigma - p)^* \langle p \rangle E$  is unambiguous.*

**Proof:** Because  $(\Sigma - p)^* \langle p \rangle (\Sigma - p)^* = \phi$ , we know that

$$(\Sigma - p)^* \langle p \rangle (\Sigma - p)^* \cap E /_{(p, E)} = \phi$$

for any  $E$ . Therefore,  $(\Sigma - p)^* \langle p \rangle E$  is unambiguous by Proposition 5.4.  $\square$

PROPOSITION 5.9. *For any regular expression  $E$  over a finite alphabet  $\Sigma$ , the extraction expression  $(\Sigma - p)^* \langle p \rangle E$  is maximal iff  $L(E) = \Sigma^*$ .*

**Proof:** First  $(\Sigma - p)^* \langle p \rangle \Sigma^*$  is unambiguous by Lemma 5.8. Because

$$((\Sigma - p)^* \langle p \rangle) \setminus ((\Sigma - p)^* \cdot p \cdot \Sigma^*) = \Sigma^*$$

and

$$((\Sigma - p)^* \cdot p \cdot \Sigma^*) /_{(p, \Sigma^*)} = (\Sigma - p)^* \cup ((\Sigma - p)^* \cdot p \cdot \Sigma^*) = \Sigma^*$$

it follows from Corollary 5.6 that  $(\Sigma - p)^* \langle p \rangle \Sigma^*$  is maximal. Thus, again by Lemma 5.8, we conclude that  $(\Sigma - p)^* \langle p \rangle E$  is maximal iff  $L(E) = \Sigma^*$ .  $\square$

THEOREM 5.10. (**Complexity of Testing Maximality**) *For any extraction expression  $E_1 \langle p \rangle E_2$  over a finite alphabet  $\Sigma$ , the problem of testing whether  $E_1 \langle p \rangle E_2$  is maximal is PSPACE-complete.*

**Proof:** From Lemma 5.7 and Proposition 5.9, we know that the problem is PSPACE-hard. Testing whether  $E_1 \langle p \rangle E_2$  is unambiguous only takes polynomial time according to Theorem 5.5. Then by Corollary 5.6, to test if  $E_1 \langle p \rangle E_2$  is maximal it suffices to test whether  $(E_1 \cdot p \cdot E_2) /_{(p, E_2)} = \Sigma^*$  and  $(E_1 \cdot p) \setminus (E_1 \cdot p \cdot E_2) = \Sigma^*$ . According to Lemma 5.2 both  $(E_1 \cdot p \cdot E_2) /_{(p, E_2)}$  and  $(E_1 \cdot p) \setminus (E_1 \cdot p \cdot E_2)$  can be computed in polynomial time. After applying Lemma 5.7 again, we conclude that the problem of testing maximality is in PSPACE.  $\square$

## 6. SYNTHESIZING MAXIMAL EXTRACTION EXPRESSIONS

In this section we first propose an algorithm, called *left-filtering maximization*, that can maximize a large class of unambiguous extraction expressions. Then we develop a *pivot maximization framework*, which can be used to enhance maximization algorithms. In particular, when applied to the left-filtering maximization algorithm, it yields a much more powerful maximization technique.

### Left-Filtering Maximization.

Consider an extraction expression  $E_1 \langle p \rangle E_2$  over a finite alphabet  $\Sigma$ . If  $(E_1 \cdot p) \setminus E_1 = \phi$ , then by Proposition 5.4 we can replace  $E_2$  with  $\Sigma^*$  and obtain a more general unambiguous extraction expression:  $E_1 \langle p \rangle E_2 \preceq E_1 \langle p \rangle \Sigma^*$ . (Similarly, if  $E_2 /_{(p, E_2)} = \phi$  then we can generalize  $E_1 \langle p \rangle E_2$  to  $\Sigma^* \langle p \rangle E_2$ .) The problem is that  $E_1 \langle p \rangle \Sigma^*$  might not be maximal, so we must do some work on  $E_1$  to make our expression maximal. The algorithm, below, is one way to maximize such an extraction expression.

#### Definition 6.1. (Finite Sequence Filtering Operator)

Given a regular expression  $E$  over a finite alphabet  $\Sigma$ , a symbol  $p \in \Sigma$ , and an integer  $n \geq 0$ , the finite sequence filtering operator  $E \|_n^p$  is defined as follows:

$$E \|_n^p \stackrel{\text{def}}{=} E \cap (\Sigma - p)^* \cdot \underbrace{p \cdot (\Sigma - p)^* \cdots p \cdot (\Sigma - p)^*}_n$$

where the suffix  $p \cdot (\Sigma - p)^*$  repeats  $n$  times.  $\square$

Informally,  $E \|_n^p$  consists of exactly those strings recognized by  $E$  that contain precisely  $n$  occurrences of the symbol  $p$ . Since the intersection of two regular expressions can be computed in polynomial time,  $E \|_n^p$  can be computed in polynomial time.

#### ALGORITHM 6.2. (Left-Filtering Maximization)

**Input:** an extraction expression  $E \langle p \rangle \Sigma^*$ , where  $\Sigma$  is a finite alphabet and  $E /_{(p, \Sigma^*)} \|_n^p = \phi$ , for some  $n \geq 0$ .  
**Output:** a maximal, unambiguous extraction expression  $E' \langle p \rangle \Sigma^*$  that generalizes  $E \langle p \rangle \Sigma^*$ .

```

BEGIN
1   $F := E /_{(p, \Sigma^*)}$ ;
2   $S := (\Sigma - p)^* - (F \|_0^p)$ ;
3   $n := 0$ ;
4  while  $F \|_n^p \neq \phi$  {
5       $S := S + (F \|_n^p \cdot p \cdot (\Sigma - p)^* - F \|_{n+1}^p)$ ;
6       $n := n + 1$ ;
7  }
8   $E' := E + S$ ;
9  return  $E' \langle p \rangle \Sigma^*$ 
END
```

PROPOSITION 6.3. (**Correctness of Left-Filtering Maximization Algorithm**) *Let  $E \langle p \rangle \Sigma^*$  be an unambiguous extraction expression over a finite alphabet  $\Sigma$ , such that  $E /_{(p, \Sigma^*)} \|_n^p = \phi$  for some  $n \geq 0$ . Then the extraction expression  $E' \langle p \rangle \Sigma^*$  computed by Algorithm 6.2 is a maximal and unambiguous generalization of  $E \langle p \rangle \Sigma^*$ .*

### Pivot Maximization Framework.

Consider an extraction expression  $E \langle p \rangle \Sigma^*$  and suppose we can find an equivalent representation for  $E$  of the form:

$$E_1 \cdot q_1 \cdot E_2 \cdot q_2 \cdots E_n \cdot q_n \cdot E_{n+1} \quad (1)$$

such that

- $E_1\langle q_1\rangle\Sigma^*, \dots, E_n\langle q_n\rangle\Sigma^*, E_{n+1}\langle p\rangle\Sigma^*$  are all unambiguous extraction expressions; and
- $E_1\langle q_1\rangle\Sigma^*, \dots, E_n\langle q_n\rangle\Sigma^*, E_{n+1}\langle p\rangle\Sigma^*$  can be maximized to  $E'_1\langle q_1\rangle\Sigma^*, \dots, E'_n\langle q_n\rangle\Sigma^*, E'_{n+1}\langle p\rangle\Sigma^*$ , respectively.

In such a case, we shall call each  $q_i$  a *pivot* and say that  $E\langle p\rangle\Sigma^*$  is *pivot-maximizable*. It turns out (by Proposition 6.6) that given a pivot-maximizable expression as above, the expression

$$(E'_1 \cdot q_1 \cdot E'_2 \cdot q_2 \cdots E'_n \cdot q_n \cdot E'_{n+1})\langle p\rangle\Sigma^* \quad (2)$$

is a maximal and unambiguous generalization of  $E\langle p\rangle\Sigma^*$ .

Pivot maximization is a powerful framework for harnessing the various specialized maximization algorithms. For instance, the left-filtering maximization algorithm can be used in this framework as follows. Suppose  $E$  can be represented as Expression (1), where  $E_1$  matches only a bounded number of  $q_1$ 's (*i.e.*,  $E_1 /_{(q_1)\Sigma^*} \parallel_n^{q_1} = \phi$  for some  $n \geq 0$ ),  $E_2$  matches only a bounded number of  $q_2$ 's, etc., and  $E_{n+1}$  matches only a bounded number of  $p$ 's. Then conditions of left-filtering maximization apply to each of the  $E_i$ 's, so we can maximize the corresponding extraction expressions using that method. By Proposition 6.6, Expression (2) is a maximal generalization of the original expression. Note that this technique is strictly more powerful than plain left-filtering maximization. Indeed, to maximize  $E\langle p\rangle\Sigma^*$  through left-filtering, it must be the case that  $E$  matches only a bounded number of  $p$ 's. In contrast, pivot maximization requires that *only*  $E_{n+1}$  must match a bounded number of  $p$ 's, so  $E$  itself can potentially match an unbounded number of  $p$ 's.

The proof of correctness of pivot maximization relies on Proposition 6.5 (and indirectly on Proposition 6.4).

**PROPOSITION 6.4. (Composition of Unambiguous Extraction Expressions)** *If  $E_1\langle q\rangle\Sigma^*$  and  $E_2\langle p\rangle\Sigma^*$  are unambiguous extraction expressions over a finite alphabet  $\Sigma$ , where  $q, p \in \Sigma$  (and  $q = p$  is possible), then the extraction expression  $(E_1 \cdot q \cdot E_2)\langle p\rangle\Sigma^*$  is unambiguous over  $\Sigma$ .*

**PROPOSITION 6.5. (Composition of Maximal Extraction Expressions)** *If  $E_1\langle q\rangle\Sigma^*$  and  $E_2\langle p\rangle\Sigma^*$  are maximal unambiguous extraction expressions over a finite alphabet  $\Sigma$ , where  $q, p \in \Sigma$  (and  $q = p$  is possible), then the extraction expression  $(E_1 \cdot q \cdot E_2)\langle p\rangle\Sigma^*$  is maximal and unambiguous over  $\Sigma$ .*

**THEOREM 6.6. (Correctness of Pivot Maximization)** *If an extraction expression  $E\langle p\rangle\Sigma^*$  over a finite alphabet  $\Sigma$ , represented as in Expression (1), is pivot-maximizable, then Expression (2) is a maximal and unambiguous generalization of  $E\langle p\rangle\Sigma^*$ .*

**Proof:** By induction on  $n$ , the number of pivots in  $E$ . The case of  $n = 0$  is trivial.

Suppose the claim is true for  $n = k, k \geq 0$ . Then for  $n = k + 1$ , we have:

$$E\langle p\rangle\Sigma^* = (E_1 \cdot q_1 \cdots E_k \cdot q_k \cdot E_{k+1} \cdot q_{k+1} \cdot E_{k+2})\langle p\rangle\Sigma^* \quad (3)$$

Let  $F$  be  $F = E_1 \cdot q_1 \cdots E_{k-1} \cdot q_{k-1} \cdot E_k \cdot q_k \cdot E_{k+1}$ . Since the number of pivots in  $F$  is  $k$ , by induction hypothesis we know that

$$(E'_1 \cdot q_1 \cdots E'_{k-1} \cdot q_{k-1} \cdot E'_k \cdot q_k \cdot E'_{k+1})\langle q_{k+1}\rangle\Sigma^* \quad (4)$$

is a maximal and unambiguous generalization of  $F\langle q_{k+1}\rangle\Sigma^*$ .

Because  $E'_{k+2}\langle p\rangle\Sigma^*$  is a maximal expression that generalizes  $E_{k+2}\langle p\rangle\Sigma^*$  (due to the assumption that  $E$  is pivot-maximizable), Proposition 6.5 implies that the result of composing Expression (4) and  $E'_{k+2}\langle p\rangle\Sigma^*$ :

$$(E'_1 \cdot q_1 \cdot E'_2 \cdot q_2 \cdots E'_{k+1} \cdot q_{k+1} \cdot E'_{k+2})\langle p\rangle\Sigma^* \quad (5)$$

is unambiguous and maximal. Clearly, Expression (5) generalizes Expression (3), since Expression (3) is

$$(F \cdot q_{k+1} \cdot E_{k+2})\langle p\rangle\Sigma^*$$

and Expression (4) generalizes  $F\langle q_{k+1}\rangle\Sigma^*$ .  $\square$

## 7. PUTTING IT ALL TOGETHER

We are now going to revisit our motivating example of Section 3 and apply the tools and techniques we developed in Sections 4 through 6.

Consider the two HTML pages from Sections 4, in their tag-sequence representation. In both cases, we are interested in the second **INPUT**-element of the form. The corresponding extraction expressions (each one works only for one of the two pages) are as follows:

```
P H1 /H1 P FORM INPUT <INPUT> P INPUT INPUT /FORM
```

and

```
TABLE TR TD /TD TD /TD /TR TR TD /TD TD /TD /TR
FORM TR TD INPUT /TD TD <INPUT> /TD TD INPUT /TD /TR
/FORM /TABLE
```

Our strategy is to first generalize these strings into an extraction expression. Learning techniques of [18, 3, 8, 4, 5] could be utilized at this stage. For the sake of this example, we will use a simple left-to-right merging heuristic, which tries to find a sequence of tags common to the two strings and takes the union of everything in-between. This yields the extraction expression below, where we replaced `TR TD /TD TD /TD /TR TR TD /TD TD /TD /TR` with `TR ... /TR`, to save space:

$$((P H1 /H1 P) + (TABLE TR \dots /TR)) FORM (TR TD)? INPUT (/TD TD)? <INPUT> \mathcal{T}ags^* \quad (6)$$

In this expression, the symbol “?” means that the corresponding subexpression can occur zero or more times, and “+” stands for the union, as before.

By Proposition 5.4, this expression is unambiguous, but it is not maximal. If none of the heuristics succeeds in producing an unambiguous expression, then the algorithm fails. An

interesting problem is to develop heuristics for guiding the disambiguation process for extraction expressions, as mentioned in Section 8.

The left-merging heuristics used for the above example is geared towards the pivot maximization framework. In our case, we can use the symbols `FORM` and `INPUT` as pivots. It turns out that the conditions for pivot maximization are satisfied and that both of the three expressions

$$((P\ H1\ /H1\ P) + (TABLE\ TR\ \dots\ /TR))\ \langle FORM \rangle\ \mathcal{T}ags^* \\ (TR\ TD)?\ \langle INPUT \rangle\ \mathcal{T}ags^*\ (\ /TD\ TD)?\ \langle INPUT \rangle\ \mathcal{T}ags^*$$

can be maximized using the left-filtering algorithm (Algorithm 6.2). The result is the following maximal and unambiguous extraction expression:

$$(\mathcal{T}ags-\ FORM)^*\ FORM\ (\mathcal{T}ags-\ INPUT)^*\ INPUT \\ (\mathcal{T}ags-\ INPUT)^*\ \langle INPUT \rangle\ \mathcal{T}ags^*$$

It is worth noting that Expression (6) can also be maximized by a direct application of Algorithm 6.2. However, this will produce a different (much larger) extraction expression. The semantics of the two expressions will also be different: while the above expression always finds the second `INPUT`-element in the *first* form, the expression produced by a direct application of Algorithm 6.2 would be looking for a second `INPUT`-element on the page, even if the first and the second `INPUT`-element come from different forms.

## 8. CONCLUSION AND FUTURE WORK

In this paper, we have made the first few steps towards a theory of resilient data extraction from semistructured documents. To this end, we defined the notion of extraction expression, provided a correctness criterion for it (unambiguity), and formalized the intuitive notion of such an expression being “robust” in the presence of changes in the source document (maximality). We provided complexity results for deciding ambiguity and maximality and proposed powerful algorithms that can maximize very large, practical classes of unambiguous extraction expressions.

Several problems still remain. First, it is still unknown whether the general problem of maximization is decidable. Second, there still is a need for learning techniques that generate good initial unambiguous expressions that could be used by our maximization algorithms. The works discussed in Section 2 do not address this issue. One way how the existing algorithms can help the task of learning unambiguous extraction expressions is as follows: we can use them to generate ambiguous expressions first. Then we could feed this expression to a “disambiguation procedure” along with a number of counterexamples. Developing such disambiguation techniques is a topic for future research. The third line of work is to explore classes of regular expressions that can be maximized with lower computational complexity.

Another interesting issue is to explore data extraction from XML. Although XML documents are generally better structured than HTML, automatic extraction from such documents calls for the creation of ontologies in order to be able

to generate suitable queries automatically. We believe that the techniques presented here along with the works on learning regular expressions discussed in Section 2 can provide a simpler solution. However, XML can make this task simpler and more reliable. One interesting issue here is using DTDs to guide the learning algorithms.

Finally, we should mention that, like all extraction techniques that are based on regular expressions, our framework has limitations. For instance, we cannot learn or generalize extraction expressions that can be expressed only using context-free grammars. A typical example here is extracting the *middle row* from dynamically generated tables. Indeed, the training set for such a learning system would consist of extraction expressions of the form `TR < TR > TR`, `TRTR < TR > TRTR`, etc. The desired pattern to learn here is `TRn < TR > TRn`, but the language recognized by this expression is not regular, so this extraction problem cannot be solved using regular expression based techniques. It would therefore be interesting to extend our framework to include more general patterns. It would then be possible to apply our results to works like [6] and, thus, enhance their results.

## Acknowledgment.

The authors would like to thank C.R. Ramakrishnan for the helpful comments and suggestions.

## 9. REFERENCES

- [1] S. Abiteboul. Querying semi-structured data. In *Int'l Conference on Database Theory*, volume 1186, pages 1–18, Delphi, Greece, 1997. Springer.
- [2] Serge Abiteboul, Dallon Quass, Jason McHugh, Jennifer Widom, and Janet L. Wiener. The lorel query language for semistructured data. *Int. J. on Digital Libraries*, 1(1):68–88, 1997.
- [3] B. Adelberg. NoDoSe: A tool for semi-automatically extracting structured and semi-structured data from text documents. In *ACM SIGMOD Conference on Management of Data*, pages 283–294, Washington, 1998. ACM.
- [4] D. Angluin. Finding patterns common to a set of strings. In *ACM Symposium on Theory of Computing*, pages 130–141, 1979.
- [5] D. Angluin. Inductive inference of formal languages from positive data. *Information and Control*, 45:117–135, 1980.
- [6] N. Ashish and C. Knoblock. Wrapper generation for semi-structured internet sources. *ACM SIGMOD Record*, 26(4):8–15, 1997.
- [7] P. Atzeni and G. Mecca. Cut & paste. In *ACM Symposium on Principles of Database Systems*, pages 117–121, Arizona, June 1997. ACM.
- [8] R.C. Berwick and S. Pilato. Learning syntax by automata induction. *Machine Learning*, 2:9–38, 1987.

- [9] B.Ribeiro-Neto, A.H.L. Laender, and A.S. da Silva. Extracting semi-structured data through examples. In *Proceedings of the International Conference on Knowledge Management*, November 1999.
- [10] P. Buneman. Semistructured data. In *ACM Symposium on Principles of Database Systems*, pages 117–121, Tucson, Arizona, June 1997.
- [11] P. Buneman, S. Davidson, G. Hillebrand, and D. Suciu. A query language and optimization techniques for unstructured data. In *ACM SIGMOD Conference on Management of Data*, Montreal, Canada, 1996. ACM.
- [12] D.W. Embley, D.M. Campbell, Y.S. Jiang, S.W. Liddle and D.W. Lonsdale, Y.-K. Ng, and R.D. Smith. Conceptual-model-based data extraction from multiple-record web pages. *Journal of Data and Knowledge Engineering*, November 1999.
- [13] Jean-Robert Gruser, L. Raschid, M. E. Vidal, and L. Bright. Wrapper generation for web accessible data sources. In *Proceedings of the Third International Conference on Cooperative Information Systems (CoopIS98)*, pages 14–23, New York City, New York, 1998.
- [14] J. Hammer, H. Garcia-Molina, J. Cho, A. Crespo, and R. Aranha. Extracting semistructured information from the web. In *In Proceedings of the Workshop on Management of Semistructured Data*, pages 18–25, Tucson, Arizona, May 1997.
- [15] J. Hammer, Hector Garcia-Molina, S. Nestorov, Ramana Yerneni, Markus M. Breunig, and Vasilis Vassalos. Template-based wrappers in the tsimmis system. In *ACM SIGMOD Conference on Management of Data*, pages 532–535. ACM, 1997.
- [16] J.E. Hopcroft and J.D. Ullman. *Introduction to Automata Theory, Languages, and Computation*. Addison-Wesley, Reading, MA, 1979.
- [17] <http://www.jango.com>. Jango Corporation.
- [18] N. Kushmerick, D. S. Weld, and R. B. Doorenbos. Wrapper induction for information extraction. In *Int'l Joint Conference on Artificial Intelligence*, volume 1, pages 729–737, Nagoya, Japan, 1997.
- [19] H.R. Lewis and C.H. Papadimitriou. *Elements of the Theory of Computation*. Prentice Hall, Englewood Cliffs, NJ, 1981.
- [20] Seung-Jin Lim and Yiu-Kai Ng. An automated approach for retrieving hierarchical data from html table. In *Proceedings of the International Conference on Knowledge Management*, November 1999.
- [21] M. Perkowitz, R.B. Doorenbos, O. Etzioni, and D.S. Weld. Learning to understand information on the internet: An example-based approach. *Journal of Intelligent Information Systems*, 8(2):133–153, March 1997.
- [22] A. Sahuguet and F. Azavant. Web Ecology: Recycling HTML pages as XML documents using W4F. In *ACM SIGMOD Workshop on Database the Web and Databases (WebBD'99)*, pages 31–35, Philadelphia, Pennsylvania, June 1999.