# XML-based Acquisition of Fine-Granular Structured Data by Grammar Morphisms

Markus Lepper, Baltasar Trancón y Widemann, Jacob Wieland

Technische Universität Berlin, Fakultät IV, Institut für Softwaretechnik und Theoretische Informatik, ÜBB, Sekr. FR 5–13, Franklinstr. 28/29, D–10587 Berlin, E-mail: {lepper,bt,ugh}@cs.tu-berlin.de

**Abstract** The DDD language and tool is a textual front-end for comfortable authoring XML documents and acquisition of structured but flexible data, esp. suited for legacy problems.

DDD has a two-layered architecture: The upper layer contains explicit mark-up for ergonomic data input, characterized by single-character-escaping and tag inference. In contrast the lower layer does not use mark-up, but parser definitions by extended regular expressions, which implicitly construct the data model from every-days textual representation of structured entities like addresses, names, domain specific notions etc.

The algorithms for transforming the front-end grammar defining terms into abstract data definitions and then into XML DTD definitions are given, as well as the algorithm for the parser. All algorithms are based on term rewriting.

**Keywords:** Data Acquisition, Semi-Formal Data, Accessibility, *inter language working*, XML, *Document Object Model*, DTD generation, SCHEMA

## 1 Goals, Contexts and Related Work

### 1.1 Goals of the DDD Research and Tool Project

The work presented herein is part of the DDD project [8]. DDD aims at making legacy text corpora as well as new out-comings of traditional text editing and authoring techniques accessible for automated text processing. This is achieved by adding the *minimum* of necessary mark-up, thereby presenting the text as „semi-structured data" to automated processing, while at the same time keeping it readable and writable by humans.

The upper layer of DDD uses *explicit* tagging, characterized by (1) one single escape character and (2) context dependent tag inference, — two rather simple principles yielding astonishingly complex mathematical consequences, cf. [8].

```
#adressbook
  #entry  ugh  #name Wieland, Jacob #phone +49 30 / 314 -23456
               #adress TU Berlin, FR 5-30
  #entry  bt ...
```

which is translated into a structure like

```
<adressbook:entry adressbook:id="ugh">
   <name>Wieland, Jacob</name>
   <phone>+49 30 / 413 -23456 /phone>
   <adress> TU Berlin, FR 5-30</adress>
</adressbook:entry>
<adressbook:entry> ...
```

The definition of the structures are given as terms of a small declaration language (synchronously generating the inference rules for parsing *and* a data model), which then is mapped to an XML DTD in a canonical way, but behaves itself somehow more orthogonal than the latter.

So, w.r.t. *data acquisition* and derivation of *data-base* meta-models, only one single *text*-meta-model (given anyhow, as DDD schema or DTD or XML schema) is given, and there is ongoing research by others dealing with transformations, consistency conditions, derivations etc. between these two, cf. [?],[?].

## 1.2   The Parser Layer of DDD

The parser from foregoing example, using only explicit tagging, delivered the „phone number" just as an unstructured string of characters. This can hardly be called „data acquisition". On the other hand we do not want to tag such fine-granular data explicitly. It seems neither economic nor convenient to be forced to write this like

```
#phone #country 49 #city 30 #main 314 #extension 45678
```

Fine granular data in semi-formal documents is traditionally represented by relatively small text fragments which follow a certain grammar, like personal names, postal addresses, calendar dates, citations, domain specific abbreviations etc. Our proposal is, to re-model the *grammar* of the existing every-day's encoding as a formal language, — also to make legacy documents treatable with least effort[1].

The syntax of these constructs is in most cases *regular*, that is neither context-dependent nor even requiring recursive definitions, — nevertheless we provide means for recursive definition, which indeed will require most of the calculation effort.

Using standard regular expressions (plus a set of arbitrary given „Scanners") we can write a definition like ...

```
phonenumber =  ("+" #decDigit #decDigit ) (#decDigit +)
               "/" (#decDigit +)  ("-" (#decDigit +) )?
```

---

[1] These grammars are in many cases highly dependent of the *cultural* context of the data's origin. One hypothesis back-grounding the DDD project is that the re-modeling of these traditionally grown grammars as exact expressions of a formal language maybe serve as analyzation technique w.r.t. the underlying structure of the denotated data itself, which may also be culturally determined.

The resulting parser will match the phone number from the example above.

One central design goal of DDD is to avoid the necessity of redundant definitions, and concentrate all necessary declarations as much as possible. So we extend the front-end grammar defining regular expressions by a means for specifying *where to put* the matched character data for a given sub-expression.

The top-level data model is that of objects (or „Elements“), to which „some data content“ is linked to *associations*, which are indexed by names. This names are called „Tags“ in the following, — and will mostly indeed be mapped to XML „Tags“.

The type of the data delivered by applying a tag name to the association belonging to some source object is defined with the definition of this „source“ object's type.

This data type can be just another element type, i.e. a singleton set containing one of the elements of this element type, or may be a *bag* of such objects, or a *sequence*, or even a sequence containing objects of *different* types.

Following the principle of „least points of declaration“, we enhance the notation for defining the front-end syntax by a means for prescribing the *tag*, the corresponding data structure to which will receive the matched content.

In our example we can prescribe these target tags of the matched character data by writing . . .

```
phonenumber = ("+" [Country #decDigit #decDigit ])? [Area #decDigit +]
              "/" [Main #decDigit +]  ("-" [Extension  #decDigit +] )?
```

What this really means, „putting“ the parsed character data into a data object identified by the given tags, depends on the data type of the element/tag combination, the inference of which is the main topic of this paper.

## 1.3   Data Model Derivation

But now the picture w.r.t. data modeling is more complex than in the tagged layer described above, because now *three* declarations have to be synchronized: (1) The grammar of the front-end, (2) the mathematical data model, and (3) the XML encoding.

It has to be analyzed, how the explicitely given tag indications together with the implicit behavior of the grammar (which cannot be easily overlooked even by a skilled user) can be mapped to a single mathematical data model.

This model has to be as fine-granular as necessary, but also as abstract as possible. The latter, because it is mostly not feasible to store all information of the front-end representation in a data base system efficiently, — generally it is not even desirable to do so, because the meant data content shall of course be realized independently from the accidentally chosen front-end representation.

In the example above, one can infer that there are tags Country, Area, etc., which all point to one single „character data“ object each. Consider e.g. one wants to encode more than one extension, which can easily achieved by just replacing a ? by an * :

```
phonenumber = ("+" [Country #decDigit #decDigit ])? [Area #decDigit +]
              "/" [Main #decDigit +]  ("-" [Extension #decDigit +] )*
```

The mathematical model inferable from this definition is now just a *bag* of numbers, — but if the *order* of these numbers shall carry semantics, we have to derive a *sequences*. This fact must of course be indicated by the author of the data model, using a further keyword:

```
phonenumber = ("+" [Country #decDigit #decDigit ])? [Area #decDigit +]
              "/" [Main  #decDigit +]
                                (#ORD "-" [Extension  #decDigit +] )*)
```

This declaration will generate an *ordered* list of extension numbers.

More complex declarations can be written using non-terminals, or defining tag contents explicitly. Both methods can be mixed freely, as shown in . . .

```
phonenumber = country ? Area  "/" main    (#ORD extension *)

country     = "+" [Country #decDigit #decDigit ]
[Area]      = #decDigit +
main      = [Main #decDigit +]
extension    = "-" [Extension  #decDigit +]
```

## 2   Construction of Meta-Models and Transformations

### 2.1   Grammar Definition

The starting definitions of our front-end representation shall be given as a *module* $\mathcal{M}$.

First we construct a set of identifiers $T$ which will be used as „Tag" values, which may occur in the derived data model, and a set $N$ of names of „non-terminals" used for defining the front-end grammar. Further shall there be a set of predefined (or maybe user-defined) „scanners", which later will be used to consume the front-end character data, thus representing the „terminals" of our grammars[2].

All sets shall be disjoint[3], and no identifier of the former sets may start with some given „reserved" character, for which we chose the „#".

$$
\begin{array}{|l}
T : \mathbb{F}\ Ident \\
N : \mathbb{F}\ Ident \\
Scan : \mathbb{F}\ Ident \\
\hline
\mathsf{disjoint}\ \langle N, T, Scan \rangle \\
\forall\, t \in (N \cup T) \bullet t.1 \neq \texttt{\#}
\end{array}
$$

---

[2] The topic of user-defined scanner generation, and the related problems of white-space-preservation, mixed content generation etc. are not treated in this paper.

[3] We will use Z notation in the following, slightly enhanced (e.g. by type sums) and slightly weakened for sake of readability.

The set of „*Extended Regular Expressions*" is now given by the free type definition:

$$\mathcal{X} \quad ::= \mathcal{X} :: \mathcal{X} \ | \ \mathcal{X}|\mathcal{X} \ | \ \mathcal{X}* \ | \ \mathcal{X}+ \ | \ \mathcal{X}? \ | \ \#\texttt{Perm}(\mathbb{N}\mathbb{N}\mathcal{X} \ \ldots \ \mathbb{N}\mathbb{N}\mathcal{X})$$
$$| \ [T \ \mathcal{X}] \ | \ \#\texttt{ORD} \ \mathcal{X} \ | \ Scan$$

$$Scan == \#\texttt{Ident} \ | \ \#\texttt{ident} \ | \ \#\texttt{Int} \ | \ \#\texttt{Float} \ | \ "stringConstant"$$
$$| \ Enum \ | \ \ldots$$

Our „extended" RegExps firstly behave just normal regular expressions: In the first line :: means sequencing (indicated as mere juxtaposition in standard Regexp, or e.g. as a comma-operator „," in XML-DTD), | means alternative, and also *,+,? have their standard meaning[4].

Please note that our XRegExps need *not* to be „un-ambiguous" nor their languages „LL(1)". The semantics we give to XRegExps (cf. the implementation of the parser in section 3) resolve ambiguity by applying the principle of longest prefix match w.r.t. the *whole* XRegExp. For instance, when p has the following definition

```
 p = (  a a (a|b) (a|b) | a a ) ( b b b b b )?
//                      a a    b b b b    ccc
//      a a    b    b                     bbbccc
```

then matching of aabbbbbbccc against p will chose the *shorter* prefix for matching the first „()", as indicated in the first comment line, because the *overall* match will be longer.

Additionally we have three „interfaces" for controlling the data-flow, two „upward" and one „from below": The notation [ t x ] means that the content which (when parsing some text input) afterwards will be recognized as an expansion of $\alpha$, shall be assigned as content to an element with the tag value $t$.

The central feature ($\#\texttt{ORD}\,x$) is used to indicate explicitly that the *order* of the content carries some semantics.

Finally the call of a scanning function $s \in Scan$ is used to control the final parsing process of the input character data.

Then a definition module is given as a function

$$\mathcal{M} \ : \ \textsf{Module} \ = \ (T \cup N) \ \to \ \mathcal{X}$$

## 2.2 Term Transformation

For manipulating the definition module we first define a basic rewriting function for XRegexp, which is given in figure 1. This function can be used for actually rewriting a term, or just for collecting information. As given in the figure, visit

---

[4] The #Perm operator is intended to make some evaluation structures directly usable in the front-end declaration language and is not treated in this paper.

**Figure1** Visitor Function for XRegexp

$$
\begin{array}{l}
\rule{7cm}{0.4pt}\\[-2pt]
\text{\textit{RewriteXRegexp}}[\tau]\\[4pt]
\textit{visit} : \tau \times \mathcal{X} \rightarrow \tau \times \mathcal{X}\\
\textit{// should be: } \textit{visit}[\beta] \ : \ \tau \times \beta \times \mathcal{X} \ \rightarrow \ \tau \times \beta \times \mathcal{X}\\
\textit{//} \qquad\qquad \textit{visit} \ = \ \textit{visit}[\,]\\[6pt]
x : \mathcal{X}, \pi : \tau, n : N, t : T, l : (N \cup T) \nrightarrow \mathbb{N} \times \mathbb{N}\\[6pt]
\begin{aligned}
\mathsf{visit}_\pi\, t \quad &= \pi, t\\
\mathsf{visit}_\pi\, n \quad &= \pi, n\\
\mathsf{visit}_\pi\, [\, t\ x\, ] \quad &= (\mathsf{visit}_\pi\, x).1,\ \ [\, t\ (\mathsf{visit}_\pi\, x).2\ ]\\
\mathsf{visit}_\pi\, \texttt{\#ORD}\ x \quad &= (\mathsf{visit}_\pi\, x).1,\ \ \texttt{\#ORD}\ (\mathsf{visit}_\pi\, x).2\\
\mathsf{visit}_\pi\, x\ \texttt{*} \quad &= (\mathsf{visit}_\pi\, x).1,\ (\mathsf{visit}_\pi\, x).2\ \texttt{*}\\
\mathsf{visit}_\pi\, x\ \texttt{+} \quad &= (\mathsf{visit}_\pi\, x).1,\ (\mathsf{visit}_\pi\, x).2\ \texttt{+}\\
\mathsf{visit}_\pi\, x\ \texttt{?} \quad &= (\mathsf{visit}_\pi\, x).1,\ (\mathsf{visit}_\pi\, x).2\ \texttt{?}\\
\mathsf{visit}_\pi\, x_1\ x_2 \quad &= (\mathsf{visit}_{\pi'} x_2).1,\ (\mathsf{visit}_\pi\, x_1).2\ ::\ (\mathsf{visit}_{\pi'}\, x_2).2\\
\mathsf{visit}_\pi\, x_1\ \texttt{|}\ x_2 \quad &= (\mathsf{visit}_{\pi'} x_2).1,\ (\mathsf{visit}_\pi\, x_1).2\ \texttt{|}\ (\mathsf{visit}_{\pi'}\, x_2).2\\
\mathbf{where}\ \pi' \ &== \ (\mathsf{visit}_\pi\, x_1).1
\end{aligned}
\end{array}
$$

is just a complicated way of writing the *identity* function. Derivations from this visitor are made by instantiating it with the types of the accumulator parameters, and overriding the function definitions, where all non-overridden methods are „inherited“ from figure 1, e.g. the cases of the function definition are treated like „methods“ in an OO-context[5].

E.g. the function „flat“ collects all values occurring at the leaves of an XRegexp and shows the typical and convenient way of working with „functionally encapsulated pseudo side-effects“:

$$
\begin{aligned}
\mathsf{flat}\ &:\ \mathcal{X}\ \rightarrow\ \mathbb{P}(N \cup T)\\
\mathsf{flat}\,(x)\ &=\ \mathsf{RewriteXRegexp}\,[\mathbb{P}(N \cup T)].\mathsf{visit}\\
&\qquad \bigoplus\ \big(\mathsf{visit}_S(t : T)\ =\ S \cup \{t\}\\
&\qquad\qquad\ \mathsf{visit}_S(n : N)\ =\ S \cup \{n\}\\
&\qquad\ )\ (\{\}, x)\ .1
\end{aligned}
$$

### 2.3 Normalization of Order-Respecting Directives

First step of the transformation of a given module $m \in \mathcal{M}$ is to remove all `#ORD` constructors, keeping in each branch only the top-level ones. This is done

---

[5] For convenience we only give one single definition of a multi-purpose visitor. When either the accumulator parameters or the rewriting facility is not used in the following, we simply omit the corresponding parameters and results from the notation informally.

by applying normOrder to all those $x \in \mathcal{X}$ which appear on the right side of definitions in $m$.

Since the #ORD constructor creates an order-respecting area which shall distribute even over the *contents* of all enclosed non-terminals, we capture this fact by doubling the set of non-terminals and their definitions, to be used in ordered and un-ordered contexts resp.

$$
\begin{aligned}
N_{\mathsf{o}} &= \{\, n \in N \bullet n_{\mathsf{o}} \,\} \\
\mathcal{M}_1 &= \mathcal{M} \cup \{\, n \mapsto x \in \mathcal{M} \bullet n_{\mathsf{o}} \mapsto \mathsf{elimOrd}(x) \,\} \\
\mathsf{normOrder}, \mathsf{elimOrder} &: \ \mathcal{X} \rightarrow \mathcal{X} \\
\mathsf{normOrder} &= \mathsf{RewriteXRegexp}\,[].\mathsf{visit} \\
&\oplus \big( \mathsf{visit}\,(\texttt{\#ORD}\,x) = \texttt{\#ORD}(\mathsf{elimOrder}(x)) \\
&\qquad )
\end{aligned}
$$

$$
\begin{aligned}
\mathsf{elimOrder} &= \mathsf{RewriteXRegexp}\,[].\mathsf{visit} \\
&\oplus \big( \mathsf{visit}\,\texttt{\#ORD}\,x = x \\
&\qquad \mathsf{visit}\,(n : N) = n_{\mathsf{o}} \\
&\qquad \mathsf{visit}\,[\,t\ \alpha\,] \ \mid\ \alpha \notin \mathsf{ran}\,\texttt{\#ORD} = [\,t\ (\texttt{\#ORD}\,\mathsf{visit}(\alpha))\,] \\
&\qquad )
\end{aligned}
$$

From here on we simply write $N$, meaning $N \cup N_{\mathsf{o}}$.

### 2.4 Extraction of In-Line Content Models

Now, as the „ordered" information is pushed down into the XRegExp, we can extract and collect all definitions for element content, which are written „in-line" inside an XRegExp:

$$
\begin{aligned}
\mathsf{exCD1} &= \mathsf{RewriteXRegexp}\,[T \times \mathcal{X}].\mathsf{visit} \\
&\oplus \big( \mathsf{visit}_{\bot,\bot}\,[\,t\ \alpha\,] = t, \alpha, \texttt{CALL}(t) \\
&\qquad \mathsf{visit}_{t,\alpha}\,x = t, \alpha, x \\
&\qquad )
\end{aligned}
$$

$$
\mathsf{exCD2} : \mathsf{Module} \rightarrow \mathsf{Module}
$$

$$
\mathsf{exCD2}(M) = \left\{
\begin{array}{l}
M \oplus (t_i \mapsto c_i) \oplus (t_o \mapsto c_o) \\
\qquad \textbf{if } \exists\, t_o \in (N \cup T) \\
\qquad\qquad\qquad \bullet\ \mathsf{exCD1}(\bot, \bot, M[t_o]) = (t_i, c_i, c_o) \wedge t_i \neq \bot \\
M \ \textbf{otherwise}
\end{array}
\right\}
$$

$$
\mathsf{extractContentDefs} : \mathsf{Module} \rightarrow \mathsf{Module}
$$
$$
\mathsf{extractContentDefs} = \mathbf{Y}\ \mathsf{exCD2}
$$

From here on we assume that $\mathcal{M}_2 = \mathsf{extractContentDefs}(\mathcal{M}_1)$, and simply write $\mathcal{M}$ for $\mathcal{M}_2$.

### 2.4.1 Example With the last both steps a definition module like . . .

```
[C] = b  b  [C #ORD b  (#int | (#ORD b*) ) ]
b   = #int [C #int] | (#ORD #ident *)
```

. . . is transformed to . . .

```
[C] = b  b  [C #ORD bo  (#int | (      b*) ) ]
b   = #int [C #int] | (#ORD #ident *)
bo  = #int [C #int] | (       #ident *)
```

. . . and then the content definition for C are extracted .

```
[C] = b  b  C
    | #ORD b  (#int | (      bo*) )
    | #int
b   = #int C | (#ORD #ident *)
bo  = #int C | (       #ident *)
```

## 2.5   Abstracting Grammars to Cardinalities

The next steps do the abstraction from regular expressions to just counting the possible occurrences of Tags and Nonterminals: What we do is building a kind of „disjunctive normal forms": For each given XRegExp (which is bound to the declaration of an tag or nonterminal) we consider all possible „paths".

Since only net „effects" are considered, and order must be discharged, we switch now to modeling these normal forms as „sets"[6].

The algorithm is given in figures 2 and 3. We model every „path" through an XRegExp, i.e. every term of our virtual disjunctive normal form, as a pair of two sets, representing the ordered and unordered occurrences, resp., each given the minimal and the maximal count.

In a first step the constructors ?, + and * are resolved. Then all paths in which recursion occurs are lifted to „infinity", and all XRegExp where a „power recursion" occurs, (or at least two distinct simple recursions in one same path!), are lifted to „infinity" for their *whole* content.

Then all nonterminals are replaced by there summarized net content recursively, and at last all paths for a given Nonterminal/Tag are summarized.

## 2.6   Possible Data-Base Content Models

Now we have indeed a function

$$\mathsf{SumFin} \ : \ (T \cup N) \ \rightarrow \ (\ ((T \cup Parse) \ \rightarrow \ (\mathbb{N} \times \mathbb{N})) \times ((T \cup Parse) \ \rightarrow \ (\mathbb{N} \times \mathbb{N}))\ )$$

giving for any Tag value ($t \in T$) or any Non-Terminal ($n \in N$) the minimal and maximal multiplicity, in which a given parser or tag reference may occur in valid contents of $n$ or $t$. This is given for „ordered" and „unordered" occurrences separately.

---

[6] Please note that using „sets" is just a simplified way of writing term transformations, e.g. modulo a certain permutation order. All data („sets" and „relations") are given by construction, and all evaluation can always be done symbolically by transformations.

**Figure2** Collecting net effects per path

---

effect $\qquad = T \nrightarrow (\mathbb{N} \times \mathbb{N})$

effects $\qquad =$ effect $\times$ effect

*// repr. ORDered and UNORDered .*

$combine \qquad :$ effects $\times$ effects $\times$ (effects $\rightarrow$ effects)$^2 \;\rightarrow\;$ effects

$combine(e_1, e_2, \phi_1, \phi_2)$

$\qquad =$ dom $e_2 \lhd e_1 \;\cup\;$ dom $e_1 \lhd e_2$

$\qquad\quad \cup \, (\lambda\, x \bullet x \mapsto ((\phi_1(e_1(x).1, e_2(x).1), \phi_2(e_1(x).2, e_2(x).2))$

$\qquad\qquad (\!|\, \text{dom } e_1 \cap \text{dom } e_2 \,|\!)$

$\uplus, \diamond \qquad\qquad :$ effects $\times$ effects $\rightarrow$ effects

$\uparrow, \downarrow \qquad\qquad :$ effects $\rightarrow$ effects

$\hat{*}, \downarrow \qquad\qquad :$ effects $\times \mathbb{N} \times \mathbb{N} \rightarrow$ effects

$e_1 \uplus e_2 \qquad = combine(e_1, e_2, (+), (+))$

$e_1 \diamond e_2 \qquad = combine(e_1, e_2, (min), (max))$

$e \uparrow \qquad\quad = (\lambda\, x \bullet x \mapsto (e(x).1, \infty)) \;\; (\!|\, \text{dom } e \,|\!)$

$e \downarrow \qquad\quad = (\lambda\, x \bullet x \mapsto (0, e(x).2)) \;\; (\!|\, \text{dom } e \,|\!)$

$e\hat{*}(a, b) \qquad = (\lambda\, x \bullet x \mapsto (e(x).1 * a, e(x).2 * b)) \;\; (\!|\, \text{dom } e \,|\!)$


effs, recEffs, finEffs $\;: \mathcal{X} \rightarrow \mathbb{P}$ effects

effs $(n : N) \qquad = \{(\{\}, \{n \mapsto (1,1)\})\}$

effs $(t : T) \qquad = \{(\{\}, \{t \mapsto (1,1)\})\}$

effs $(x \,|\, y) \qquad = \lambda\, a, b \bullet ((a.1 \diamond b.1), (a.2 \diamond b.2)) \;(\!|\,\text{effs } x \times \text{effs } y \,|\!)$

effs $(x :: y) \qquad = \lambda\, a, b \bullet ((a.1 \uplus b.1), (a.2 \uplus b.2)) \;(\!|\,\text{effs } x \times \text{effs } y \,|\!)$

effs $(x?) \qquad = \_ \downarrow \;\; (\!|\,\text{effs } x \,|\!)$

effs $(x{+}) \qquad = \_ \uparrow \;\; (\!|\,\text{effs } x \,|\!)$

effs $(x*) \qquad = \_ \downarrow \; (\!|\; \_ \uparrow \; (\!|\,\text{effs } x \,|\!) \,|\!) \;\; = \;\; \text{effs}((x{+})?)$

effs $(\texttt{\#ORD}\,x) \qquad = (\lambda\, s \bullet (s.2, \{\})) \;\; (\!|\,\text{effs } x \,|\!)$

*// bec. of normalization there is only UNORDered content below*

*// — shift it to ordered content !*


$usesMulti, usesOnce, uses \;: $ effect $\leftrightarrow N$

$usesMulti \qquad = \; \{(e, t) \;|\; (e.1)(b).2 \,>\, 1 \;\; \vee \;\; (e.2)(b).2 \,>\, 1 \}$

$usesOnce \qquad = \; \{(e, t) \;|\; (e.1)(b).2 \,=\, 1 \;\; \vee \;\; (e.2)(b).2 \,=\, 1$

$\qquad\qquad\qquad\qquad \wedge \; \neg \; usesMulti(a, b)\}$

$uses \qquad\qquad = usesOnce \;\cup\; usesMulti$


$calls, callsDirectMulti, callsDirectOnce, inCircleWith$

$\qquad\qquad\qquad : \quad N \leftrightarrow N$

$callsDirect \qquad = \; \{(a, b) \;|\; \exists\, e \in \text{effs } \mathcal{M}[\text{a}] \bullet uses(e, b)\}$

$calls \qquad\qquad = (callsDirect)^*$

$callsDirectMultiple = \; \{(a, b) \;|\; \exists\, e \in \text{effs } \mathcal{M}[\text{a}]$

$\qquad\qquad\qquad\quad \bullet usesMulti(e, b)$

$\qquad\qquad\qquad\quad \vee \; (uses(e, b) \wedge \exists\, m \in N \,|\, m \neq b \bullet e(m).1.2 + e(m).2.2 > 0 \wedge \; m\; calls\; a)$

$\qquad\qquad\qquad \}$


$callsDirectSingle \;=\; callsDirect \,\setminus\, callsDirectMultiple$

$callsMultiple \;=\; \{(a, d) \;|\; \exists\, b, c : N \bullet \qquad (a\; calls\; b \;\vee\; a = b)$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad \wedge (b\; callsDirectMultiple\; c)$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad \wedge (c\; calls\; d \;\vee\; c = d)$


isMultiRec $\qquad : \mathbb{P}\, N$

$\qquad\qquad\qquad = \{n : N \;|\; n\; callsMultiple\; n \}$

---

**Figure3** Collecting net effects per path *(continued)*

$$\mathsf{recEffs}(n) \quad = \left\{ \begin{array}{ll} (\mathsf{effs}\ n)\uparrow & \textbf{if } \mathsf{isMultiRec}(n) \\ (\lambda\, x \bullet \mathit{calcRec}(n,x,\mathcal{M}))\ (\!|\ \mathsf{effs}(\mathcal{M}[n])\ |\!) & \textbf{otherwise} \end{array} \right\}$$

$$\mathit{calcRec} \quad\quad\quad : N \times \mathsf{effect} \times \mathsf{Module}\ \rightarrow \mathsf{effect}$$

$$\mathit{calcRec}(n,e,\mathcal{M}) = \left\{ \begin{array}{ll} e\uparrow & \textbf{if } \exists\ n' \in \mathcal{M} \bullet e\ \mathit{uses}\ n'\ \wedge\ n'\ \mathit{calls}\ n \\ e & \textbf{otherwise} \end{array} \right\}$$

$$\mathsf{sumEff},\mathsf{sumFin} \quad : N\ \rightarrow\ \mathsf{effect}$$

$$\mathsf{sumEff}(n) \quad\quad = \mathit{collEff}(\mathsf{recEffs}\ \mathcal{M}[n])$$

*// accumulate effect of the alternatives !*

$$\mathit{collEff} \quad\quad\quad\quad : \mathsf{effects}\ \rightarrow\ \mathsf{effect}$$

$$\mathit{collEff}\ (\{e_0\} \cup e) = e_0\ \diamond\ \mathit{collEffs}(e)$$

$$\mathit{collEff}\ (\{\}) \quad\quad = \quad (\{\},\{\})$$

$$\mathsf{sumFin}(n) \quad\quad = \mathbf{Y}\ \mathit{collFin}(\{n\},\mathsf{sumEff}\ \mathcal{M}[n])$$

*// replace nonterminals by terminal content.*

$$\mathit{collFin} \quad\quad\quad\quad : \mathsf{effect} \times \mathbb{P}\, N\ \rightarrow\ \mathsf{effect}$$

$$\mathit{collFin}(e,S) \quad\quad = \left\{ \begin{array}{l} (\{m\} \lhd e) \uplus\ (\ \mathit{collFin}(\mathcal{M}[m], S \cup \{m\})\ \hat{*}(a.b)) \\ \quad\quad \textbf{if}\, \exists\, m : N \bullet m \mapsto (a,b)\ \in e\ \wedge\ m \notin S \\ e\ \textbf{otherwise} \end{array} \right\}$$

From this data we can derive the abstract definitions of the data model, which is appropriate for storing the parsed data according to the user's prescriptions:

| $\mathcal{O}$ | | $\mathcal{U}$ | | | |
|---|---|---|---|---|---|
| 0 | 0 | 1 | 1 | $b:\ \Delta \rightarrow \mathcal{T}$ | $—\!\!\diamond^1\!\!—\mathcal{T}$  **maybeXML − ATT** |
| 0 | 0 | 0 | 1 | $n:\ \Delta \rightarrow \mathtt{OPT}\ \mathcal{T}$ | $—\!\!\diamond^{0..1}\!\!—\mathcal{T}$**maybeXML − ATT** |
| 0 | 0 | $a\ b>1$ | | $n:\ \Delta \rightarrow \mathtt{BAG}\ \mathcal{T}$ | $\mathbb{N}—\!\!\diamond^1\!\!—\mathcal{T}$ |
| _ $=1$ | | _ | 1 | | |
| _ $>1$ | | _ | 1 | *// May be discussed.* | |
| _ $a>1$ | | _ | $b>1$ | **Error** *(„Automated Derivation Fails")* | |
| _ $\geq 1$ | | 0 | 0 | | |
| | | | | **let** $\mathit{ordererSubs}\ =\ \mathsf{dom}\ (\mathsf{sumFin(n)}.1)$ | |
| | | | | $\quad m\ =\ \sharp\ \mathit{ordererSubs}$ | |
| $m=1$ | | | | $n:\ \Delta \rightarrow \mathtt{SEQ}\ \mathcal{T}$ | $\mathbb{N}—^1\!\!\diamond\!\!—\mathcal{T}$ |
| $m>1$ | | | | $n:\ \Delta \rightarrow \mathtt{SEQ}\,(\mathtt{UNION}(\mathit{orderedSubs}))$ | $\mathbb{N}—^1\!\!\diamond\!\!—(\mathbb{N} \times (\mathcal{T}_1 \uplus \mathcal{T}_2 \uplus \ldots))$ |

The situation that there is one single „unordered" and more than one „ordered" occurences of the same sub-grammar may be very common. Consider e.g. the following parser, representing two possible ways of writing *personal names* taken from the cultural sphere of German language:

```
PersName = [Vorname #Ident][Nachname #ORD #Ident (#ident|#Ident)*]
         | [Nachname #ORD #Ident (#ident|#Ident)*]
               "," [Vorname #ORD #Ident +]
```

This imposes no real problems, as a *single* occurence can alwas be considered to be „ordered".

## 2.7 XML Content Models

For automated transformation we define an XML DTD to be represented as a term of the grammar given by $\mathsf{module}_X$:

$$\mathsf{module}_X ::= XmlDef \mid \mathsf{module}_X :: XmlDef$$
$$XmlDef ::= \mathtt{EL}_X(T, Xgr) \mid \mathtt{ATT}_X(T, Ident, \{\mathtt{REQUIRED}, \mathtt{IMPLIED}\})$$
$$Xgr ::= T \mid Xgr\text{ ","}\ Xgr \mid Xgr \mid Xgr \mid Xgr\ ? \mid Xgr\ + \mid Xgr\ * \mid \varepsilon_X$$

$\varepsilon_X$ is needed for technical reasons and is neutral w.r.t. sequencing, e.g.

$$x","\varepsilon_X = \varepsilon_X","x = x$$

. The transformation from $\mathsf{effects}$ to $\mathsf{module}_X$ is given in figure 4.

# 3 The Parsing Process

Parsing is realized as parallel breadth-first match, implemented again by transformation of terms, or set of terms, resp.:

The following formulae expect some globally given character data, „visible" only on Scanner level. Requests to and results from such scanners are represented by only the „position information" $p \in \mathbb{Z}$.

The function $\mathsf{parseResult}$, when called with a XRegExp from our normalized module (together with the starting position), will transform the definition term into the solution term. This term is element of the language given by

$$matched ::= [t\ matched] \mid S_{\mathbb{N} \times \mathbb{N}} \mid matched :: matched$$

it contains only references to scanners for which matching character data has been found (which is annotated as two position values), and tag constructors, which group the contained data as representing one element.

This structure can one-to-one be translated to an XML document, and fulfills all properties of the mathematical data model derived before.

## 3.1 A Simple Meta-Model for XML

Now we can give an exact definition of a well-formed XML *term*, (which can represent a well-formed document or fragment)[4]. Our definition has the additional property of disjointness of the both sets of names used for attributes („Attribute name") and for elements („Tags") resp.. We indicate this by choosing the name „**x** *Wellformed*"

**Figure4** Deriving XML DTD content models

$$[[\_]]^X \qquad\qquad\qquad : \mathsf{module} \ \to\ \mathsf{module}_X$$

$$collXml \qquad\qquad\qquad : \mathbb{P}\ T\ \to\ \mathsf{module}_X$$

$$collXD \qquad\qquad\qquad : \ T\ \to\ \mathsf{module}_X$$

$$collCounter \qquad\qquad : \ T\ \times \mathbb{N} \times \mathbb{N}_\infty\ \to\ \mathsf{module}_X$$

$$[[\,\mathcal{M}\,]]^X \qquad\qquad\qquad = collXml(\mathsf{dom}\,\mathcal{M}\ \cap\ T)$$

$$collXml\,(t\ \cup\ t_R) \qquad\qquad = collXD(t, \mathsf{sumFin}(n), \varepsilon_X, \varepsilon_X, \varepsilon_X)\ ::\ collXml(t_R)$$

$$collXml\,(\{\}) \qquad\qquad\qquad = \varepsilon_X$$

$$collXD(t, (S_O, \{\}, \{\}, \alpha, \beta, \gamma)$$
$$\qquad\qquad\qquad\qquad = \mathtt{EL}_X(t, (\beta\ "\text{,}"\ (\gamma)\,\texttt{*}\,)\,)\ ::\ \alpha$$

$$collXD(t, (S_O, \{a \mapsto (1,1)\} \cup S_U, \alpha, \beta, \gamma)$$
$$= \left\{ \begin{array}{l} collXD(S_O, S_U, \alpha :: \mathtt{ATT}_X(n, a, \mathtt{REQUIRED}), \beta, \gamma) \\ \qquad \mathbf{if}(\mathsf{mayBeAtt}\,a) \wedge (\mathsf{userWantsAtt}\,a) \\ collXD(S_O, S_U, \alpha, \beta\ "\text{,}"\ a, \gamma) \\ \qquad\qquad \mathbf{otherwise} \end{array} \right\}$$

$$collXD(t, (S_O, \{a \mapsto (0,1)\} \cup S_U, \alpha, \beta, \gamma)$$
$$= \left\{ \begin{array}{l} collXD(S_O, S_U, \alpha :: \mathtt{ATT}_X(n, a, \mathtt{IMPLIED}), \beta, \gamma) \\ \qquad \mathbf{if}(\mathsf{mayBeAtt}\,a) \wedge (\mathsf{userWantsAtt}\,a) \\ collXD(S_O, S_U, \alpha, \beta\ "\text{,}"\ (a\texttt{?}), \gamma) \\ \qquad\qquad \mathbf{otherwise} \end{array} \right\}$$

$$collXD(t, (S_O, \{n \mapsto (a, b > 1)\} \cup S_U, \alpha, \beta, \gamma)$$
$$\qquad\qquad\qquad = collXD(t, S_O, S_U, \alpha, \beta "\text{,}"\ (collCounter(n, a, b)), \gamma\,)$$

$$collXD(t, (\{a \mapsto (\_, n > 1)\} \cup S_O, S_U, \alpha, \beta, \gamma)$$
$$\qquad\qquad\qquad = collXD(t, S_O, S_U, \alpha, \beta, (\gamma "\text{,}"\ a)\,)$$

$$collCounter\,(n, a > 0, b)\ = n"\text{,}"\,(collCounter(n, a - 1, b - 1))$$

$$collCounter\,(n, 0, b \neq \infty) = (\ n"\text{,}"\,(collCounter(n, 0, b - 1)\ )\texttt{?}$$

$$collCounter\,(n, 0, 0)\qquad = \varepsilon_X$$

$$collCounter\,(n, 0, \infty)\qquad = n\ \texttt{*}$$

---

$$\boxed{\begin{array}{l} \underline{\ Xml\_xWellformed\_Term[T_E, T_A, C]\ } \\[4pt] \texttt{\#self}\ \stackrel{def}{=}\ Xml\_xWellformed\_Term[T_E, T_A, C] \\ tag\ :\ T_E \\ atts\ :\ T_A \nrightarrow C \\ cont\ :\ \mathsf{seq}\,(\texttt{\#self} \cup C) \\ \hline T_E \cup T_A \subset T \\ T_E \cap T_A = \varnothing \\ \text{.-.-.-.-.-.-.-.-.-.-.-.-} \\ sups\ :\ \mathbb{F}\ \texttt{\#self} \\ \quad = \{e : \texttt{\#self}\ |\ \exists\,n : \mathbb{N} \bullet e.cont.n = \theta\texttt{\#self}\,\} \\ \qquad \textit{// weaker:}\ |\ \theta\texttt{\#self} \in \mathsf{ran}\,(e.cont) \end{array}}$$

**Figure5** The Parsing Algorithm

$$\mathbb{N}_\perp \qquad\qquad = \mathbb{N} \cup \{\perp\}$$

$$\text{scan} \qquad\qquad : \; Scan \times \mathbb{N} \to \mathbb{N}_\perp \times \mathbb{N}$$

$$\hat{\mathcal{X}} \qquad\qquad ::= \mathcal{X} \mid s_{\mathbb{N}\times\mathbb{N}} \mid \varepsilon$$
$$\qquad\qquad\qquad \textbf{where } \varepsilon :: \alpha \;=\; \alpha :: \varepsilon \;=\; \alpha$$

$$\text{parse1} \qquad\qquad : \; \mathcal{X} \times \mathbb{N} \to \mathbb{P}(\hat{\mathcal{X}} \times \mathbb{N})$$

$$\text{parse1}(s : Scan) \;=\; \left\{ \begin{array}{l} \{\, s_{a,b} \,\} \textbf{ if } a \neq \perp \\ \{\} \qquad\quad \textbf{otherwise} \end{array} \right\}$$
$$\qquad\qquad\qquad \textbf{where }(a,b) = \text{scan}(s,p)$$

$$\text{parse1}(x\,|\,y), p \quad = \; \text{parse1}(x,p) \;\cup\; \text{parse1}(y,p)$$
$$\text{parse1}(x :: \alpha), p \quad = \; (\lambda\,\mu,\nu \bullet \mu :: \text{parse1}(\nu,\alpha)\,)\;(\!|\text{parse1}(x,p)|\!)$$
$$\text{parse1}(x?) \qquad = \; \text{parse1}(x,p) \;\cup\; \{(\varepsilon,p)\}$$
$$\text{parse1}(x\text{+}) \qquad = \; \text{parse1}(x :: (x\ast),p)$$
$$\text{parse1}(x\ast) \qquad = \; \text{parse1}((x\text{+})?,p)$$
$$\text{parse1}(n : N, p) \;= \; \text{parse1}(\mathcal{M}[n],p)$$
$$\text{parse1}(t : T, p) \;= \; [\,t\;\text{parse1}(\mathcal{M}[t],p)\,]$$

$$\text{maxParse}(x,p) \quad = \; \{(a,b) \in \text{parse1}(x,p) \mid \neg\;\exists(a',b') \in \text{parse1} \bullet b' > b \;\}$$
$$\text{parseResult}(x,p) = \; \texttt{RANDOMCHOICE}\;\text{maxParse}(x,p)$$

Standard XML only deals with finite terms[7] , given like ...

---

[7] By the way, a rather different thing is an XML document „*object*" representation, which can be described as ...

$$\underline{\qquad Xml\_xWellformed\_Object[T_E, T_A, C, \mathcal{I}]\qquad}$$

$Xml\_xWellformed\_Finite\_Term[T_E, T_A, C]$

$\texttt{\#self} \;\overset{def}{=}\; Xml\_xWellformed\_Object[T_E, T_A, C, \mathcal{I}]$

$id \;:\; \mathcal{I}$

---

$\forall\; o_1, o_2 : \texttt{\#self} \bullet o_1.id = o_2.id \;\Rightarrow\; o_1 = o_2$

$\qquad\qquad\qquad \wedge\,(o_1.cont \cap o_2.cont \neq \varnothing) \;\Rightarrow\; o_1 = o_2$

$\cdot\text{-}\cdot\text{-}\cdot\text{-}\cdot\text{-}\cdot\text{-}\cdot\text{-}\cdot\text{-}\cdot\text{-}\cdot\text{-}$

$sup : \texttt{\#self} \cup \{\texttt{\#topmost}\}$

$\qquad = \; \mu\,(\quad \{1 \mapsto \texttt{\#topmost}\,\}$

$\qquad\qquad\quad \oplus (\lambda\,e \bullet 1 \mapsto e)\;(\!|sups|\!) \;\; /\!/ = \{1\} \times subs$

$\qquad\qquad ) \,(1)$

$\theta\texttt{\#self} \;\notin\; \mathsf{ran}\,(\_.sup)^*$

markuslepper.eu

$$\boxed{\begin{array}{l} \underline{\quad Xml\_xWellformed\_Finite\_Term[T_E, T_A, C] \quad} \\ Xml\_xWellformed\_Term[T_E, T_A, C] \\ \hline // \ One\ possible\ rule\ for\ prohibiting\ circularity: \\ \theta\texttt{\#self} \notin \mathsf{ran}\ (\_.cont)^* \\ \text{.-.-.-.-.-.-.-.-.-.-.-.-.-} \\ \theta\texttt{\#self} \notin \mathsf{ran}\ (\_.sups)^* \end{array}}$$

The transformation of the data acquired by the parsing process (given in the language „*matched*") into such an XML term is straight-forward and mainly has to do the re-ordering of all „unordered" elements according to the order in the derived DTD.

## 4  Related and Future Work

### 4.1  Grammars in Data Base research

The are numerous approaches in data base research dealing with grammars and grammar transformations. This is not astonishing, if one considers grammar transformations to be an abstractions of graph transformations, which is common place in data-base theory.

The work presented herein differs (to the best of our knowledge) from all existing work, as our focus is on *front-end* generation, integrated with schema definition, and on the problems which arise when *two* independently existing grammars (front-end and data definition) interfere.

The closest relation is to the works about automated derivations of schemas (conceptual or physical) from grammars (or attributed grammars, or XML DTDs ), and vice versa, as in [15], [14], [5], [9], [7].

Only on the level of methodology there are (rather tight) relations to those approaches which use grammar morphism for transforming the *data* of semi-structured objects, as in [1] and [10], — for processing of queries, the usage of XML for presenting views and query results, and automated derivation of a corresponding DTD, as in [13], [12], — for model unification, as in [6], — and to all those approaches which propose a formal definition of SGML, XML (or semi-structured data in general) as a mathematical sound calculus, like [2] and [11].

An interesting survey on semi-structured data and XML is given in [16].

Furthermore there are common areas with *computer linguistics*, because the problem of different structural organization of front-end syntax and „data syntax" is a central topic there, cf. e.g. [3].

### 4.2  Future Work

The principle of „single point of definition", as applied in the architecture presented herein, suits well for local definitions of small grammars for „details".

It does indeed conflict with principled of *compositionality*, e.g. when using one single front-end definition with different (e.g. legacy) data definitions.

While keeping the in-lined tag definitions as an option, we need some means for combining front-end grammars with separately existing data schemas *and/or* XML schemas.

The other possible kinds of derivation (front-end grammars from DTDs or data base schemas) may also turnout to be useful.

Further basic research in these both areas seems promising.

## References

1. An Feng An and Toshiro Wakayama. SIMON: A grammar-based transformation system for structured documents. In *Fifth International Conference on Electronic Publishing*, pages 361–372, Darmstadt, Germany, 1994. Origination, Dissemination, and Design (EPODD) EP '94.

2. Chutiporn Anutariya, Vilas Wuwongse, Ekawit Nantajeewarawat, and Kiyoshi Akama. Towards a foundation for XML document databases. In *EC-Web*, pages 324–333, 2000.

3. Rens Bod and Ronald Kaplan. A data-oriented parsing model for lexical-functional grammar.

4. Tim Bray, Jean Paoli, C. M. Sperberg-McQueen, and Eve Maler. *Extensible Markup Language (XML) 1.0 (Second Edition)*. W3C Recommendation, `http://www.w3.org/TR/2000/REC-xml`.

5. Ronaldo dos Santos Mello and Carlos Alberto Heuser. A rule-based conversion of a dtd to a conceptual schema, November 2001.

6. Marc Gyssens, Jan Paredaens, and Dirk van Gucht. A grammar based approach towards unifying hierarchical data models. *SIAM Journal on Computing*, 23(6):1093–1137, 1994.

7. Gerti Kappel, Elisabeth Kapsammer, Stefan Rausch-Schott, and Werner Retschitzegger. X-ray – toward integrating xml and realational database systems. In *Conceptual Modeling – ER 2000*. Springer LNCS 1920, 2000.

8. Markus Lepper, Baltasar Trancon y Widemann, and Jacob Wieland. Minimze mark-up ! – natural writing should guide the design of textual modeling frontends. In *Conceptual Modeling — ER2001*, volume 2224 of *LNCS*. Springer, November 2001.

9. Dongwon Lee Murali Mani and Richard R. Muntz. Semantic data modeling using xml schemas, November 2001.

10. Makoto Murata. Transformation of documents and schemas by patterns and contextual conditions. In *PODP*, pages 153–169, 1996.

11. Makoto Murata, Dongwon Lee, and Murali Mani. Taxonomy of xml schema languages using formal language theory. In *Extreme Markup Languages*, `http://www.cobase.cs.ucla.edu/tech-docs/dongwon/mura0619.ps`, august 2001.

12. Frank Neven and Jan Van den Bussche. Expressiveness of structured document query languages based on attribute grammars. In *PODS '98. Proceedings of the seventeenth ACM SIGACT-SIGMOD-SIGART symposium on Principles of database systems*, pages 11–17. ACM press, 1998.

13. Yannis Papakonstantinou and Victor Vianu. DTD inference for views of XML data. In *Symposium on Principles of Database Systems*, pages 35–46, 2000.

14. Patricia Rodriguez-Gianolli and John Mylopoulos. A semantic appproach to xml-based data integration, November 2001.

15. Airi Salminen and Frank Wm. Tompa. Grammars++ for modelling information in text. *Information Systems*, 24(1):1–24, 1999.

16. Dan Suciu. Semistructured data and XML. In *FODO*, pages 0–, 1998.