

Interactive Editing of MathML Markup Using \TeX Syntax*

Luca Padovani

Department of Computer Science, University of Bologna

Mura Anteo Zamboni, 7

40127 Bologna

Italy

lpadovan@cs.unibo.it

<http://www.cs.unibo.it/~lpadovan/>

Abstract

We describe the architecture of a syntax-directed editor for authoring structured mathematical documents which can be used for the generation of MathML markup [4]. The author interacts with the editor by typing \TeX markup as he would do in a normal text editor, with the difference that the typed markup is parsed and displayed on-the-fly. We discuss issues regarding both the parsing and presentation phases and we propose implementations for them. In contrast with existing similar tools, the architecture we propose offers better compatibility with \TeX syntax, a pervasive use of standard technologies and a clearer separation of content and presentation aspects of the information.

1 Introduction

MathML [4] is an XML [2] application for the representation of mathematical expressions. Like most XML applications, MathML is unsuitable to be written directly because of its verbosity except in the simplest cases. Hence the editing of MathML documents needs the assistance of dedicated tools. As of today, such tools can be classified into two main categories:

1. WYSIWYG (What You See Is What You Get) editors that allow the author to see the formatted document on the screen while it is being composed. The editor usually provides some “export mechanism” that creates XML with embedded MathML from the internal representation of the document;
2. conversion tools that generate MathML markup from different sources, typically other markup languages for scientific documents, such as \TeX [5].

Tools in the first category are appealing, but they suffer from at least two limitations: a) editing is typically *presentation oriented* — the author is primarily concerned about the “look” of the document and tends to forget about its content. b) They may slow down the editing process because they often involve the use of menus, palettes of symbols,

and, in general, the pointing device for completing most operations.

In this paper we describe the architecture of a tool that tries to synthesize the “best of both worlds”. The basic idea is to create a WYSIWYG editor in which editing is achieved by typing concrete markup as the author would do in an actual plain text editor. The markup is then tokenized and parsed on-the-fly, a corresponding presentation is created by means of suitable transformations, and finally displayed. The editor is meant not only as an authoring tool, but more generally as an interface for math applications.

Although in the paper we assume that the concrete markup typed by the user is \TeX (more precisely the subset of \TeX concerned about mathematics) and that presentation markup is MathML, the system we are presenting is by no means tied to these languages and can be targeted to other contexts as well. One question that could arise is: “why \TeX syntax?” We can see at least three motivations: first of all because of \TeX popularity in many communities. Second, because macros, which are a fundamental concept in \TeX , are also the key to editing at a more content-oriented level, which is a primary requirement for many applications handling mathematics. Finally, because, as we will see, \TeX markup has good *locality* properties which make it suitable in the interactive environment of our concern.

* This work has been supported by the European Project IST-2001-33562 MoWGLI.

The body of the paper is structured into four main sections: in Section 2 we overview the architecture of the tool while in Sections 3, 4, 5 we describe in more detail the main phases of the editing process (lexing, parsing, and transformation). Familiarity with T_EX syntax and XML-related technologies is assumed.

2 Architecture

Several tools for the conversion of T_EX markup suffer from two major drawbacks that we are not willing to tolerate in our design: (1) they rely on the T_EX system itself for parsing the markup. While guaranteeing perfect compatibility with T_EX, this implies the installation of the whole system. Moreover, the original T_EX parser does not meet the incremental requirements that we need; (2) the lack of flexibility in the generation of the target document representation, which is either fixed by the conversion tool or it is only slightly customizable by the user.

To cope with problem (1) we need to write our own parser for T_EX markup. This is well known to be a non-trivial task, because of some fancy aspects regarding the very nature of T_EX syntax and the lack of a proper “T_EX grammar”. We will commit ourselves with a subset of T_EX syntax which appears to be just what an average author needs when writing a document. As we will see, the loss in the range of syntactic expression is compensated by a cleaner and more general transformation phase. As for the lack of a T_EX grammar, we perceive this as a feature rather than a weakness: after all T_EX is built around the fact that authors are free to define their own macros. Macros are the fundamental entities giving structure to the document.

Let us now turn our attention to problem (2): recall that the general form of a T_EX macro definition (see *The T_EXbook*, [5]) is

```
\def⟨control sequence⟩⟨parameter text⟩
  {⟨replacement text⟩}
```

where the ⟨parameter text⟩ gives the syntax for invoking the macro and its parameters whereas the ⟨replacement text⟩ defines somehow the “semantics” of the macro (typically a presentational semantics). Thus the ultimate semantic load of a macro is invariably associated with the configuration of the macro at the point of definition.

We solve problem (2) by splitting up macro definitions so that structure and semantics can be treated independently. A well-formed T_EX document can be represented as a tree whose leaves are either literals (strings of characters) or macros with no parameters, and each internal node represents a

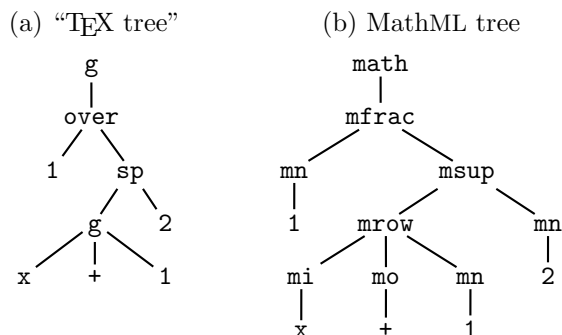


Figure 1: Tree representation for $\{1\over{x+1}\}^2$ and corresponding MathML markup.

macro and the node’s children are the macro’s parameters. Entities like delimiters, square brackets surrounding optional parameters or literals occurring in the ⟨parameter text⟩ of macro definitions are purely syntactic and need not be represented in the tree if our main concern is capturing the structure of the document. Fig. 1(a) shows the tree structure of a simple mathematical formula.

Once the document is represented as a tree, the process of macro expansion — that is, *interpretation* — can be defined as a recursive transformation on the nodes of the tree. As we will represent trees using XML, transformations can be very naturally implemented by means of XSLT stylesheets [3]. Fig. 1(b) shows the MathML tree corresponding to the T_EX tree on the left hand side. The two trees are basically isomorphic except for the name of the nodes and the presence of explicit token nodes for literals in the MathML tree. This is to say that the MathML tree can be generated from the T_EX tree by simple transformations. However, once the interpretation phase is independent of parsing (something which does not happen in T_EX) it is natural to define much more general transformations that are not just node-by-node rewritings.

The following are the main components of an interactive, syntax-based editor for structured documents:

INPUT BUFFER: the sequence of concrete characters typed by the author;

LEXICAL ANALYZER: responsible for the tokenization of the characters in the input buffer;

DICTIONARY: a map from ⟨control sequence⟩ to ⟨parameter text⟩ which is used to know the syntax of macros;

PARSER: for the creation of the internal tree structure representing the document;

TRANSFORMATION ENGINE: for mapping the internal tree into the desired format.

No doubt these entities are common to all tools converting \TeX markup into a different format, but the degree of mutual interdependence and the way they are implemented may differ considerably, especially when interactivity is a main concern. The added value of our approach is that it allows the author to independently customize both the dictionary and the transformation engine, and the advanced user of the editor the possibility of adapting the lexical analyzer to languages other than \TeX (we will spend a few more words on this topic in the conclusions).

Notation We will use the following conventions regarding lists. Lists are uniformly typed, that is elements of a list are all of the same type. We use α^* to denote the type of a list whose elements have type α . \square is the empty list; $n :: x$ is the list with head element n and tail x ; $x@y$ is the concatenation of two lists x and y ; $[n_1; \dots; n_k]$ is a short form for $n_1 :: \dots :: n_k :: \square$.

3 Lexical Analysis

The purpose of this phase is to tokenize the input buffer. As we are talking about an interactive tool, the presence of an input buffer may look surprising. Implementations for the input buffer range from *virtual* buffers (there is no buffer at all, characters are collected by the lexical analyzer which outputs tokens as they are completed) to *flat* buffers (just a string of characters as in a text editor) to *structured* buffers. For efficiency, we do not investigate in detail all the possibilities in this paper, but early experiments have shown that working with virtual buffers can be extremely difficult. As long as insert operations are performed at the right end of the buffer the restructuring operations on the parsing tree are fairly easy, but when it comes to deletion or to modifications in arbitrary positions, the complexity of restructuring operations rises rapidly to an unmanageable level. Hence, from now on we will assume that a flat input buffer is available. Whether the buffer should be visible or not is a subjective matter, and may also depend on the kind of visual feedback given by the editor on incomplete and/or incorrect typed markup.

The outcome of the lexer is a stream (list) of tokens. Each token may have one of three forms: a *literal*, that is a single character to be treated “as is”, a *space*, that is a sequence of one or more space-like characters, or a *control sequence*, that is the name of a macro.

Since the token stream is the only interface between the lexer and the parser, the lexer has the freedom to perform arbitrary mappings from the characters in the input buffer to tokens in the stream. In particular, some \TeX commands like $\backslash\alpha$ or $\backslash\rightarrow$ are just placeholders for Unicode characters. There is no point in communicating these entities as control sequences as the internal tree representation (XML) is able to accommodate Unicode characters naturally; also, treating them as literals simplifies the subsequent transformation phase.

On the other hand, there are characters, such as curly braces $\{$ and $\}$ or scripting operators $_$ and $\hat{}$, that have a special meaning. Logically these are just short names for macros that obey their own rules regarding parameters. What we propose is a general classification of parameter types which, in addition to parameters in normal \TeX definitions, allows us

- to deal with optional parameters as \LaTeX [6] does;
- to treat $\{$ as just an abbreviation for $\backslash\text{bgroup}$ and make $\backslash\text{bgroup}$ a macro with one parameter delimited by $\backslash\text{egroup}$, which we treat as the expansion for $\}$. In order for this “trick” to work we have to design the parser carefully, as we will see in Sect. 4;
- to treat scripting operators $_$ and $\hat{}$ as the two macros $\backslash\text{sb}$ and $\backslash\text{sp}$ both accepting a so-called pre-parameter (a parameter that occurs *before* the macro in the input buffer) and a so-called post-parameter (a parameter that occurs *after* the macro in the input buffer);
- to deal with macros that have “open” parameters. For instance $\backslash\text{rm}$ affects the markup following it until the first delimiter coming from an outermost macro is met. We treat $\backslash\text{rm}$ as a macro with an open post-parameter that extends as far as possible to the right. Similarly, $\backslash\text{over}$ can be seen as a macro with open pre- and post-parameters.

In order to describe parameter types we need to define the concept of *term*. A *term* is either a literal or a macro along with all its parameters (equivalently, a term is a subtree in the parsing tree). A *simple* parameter consists of one sole term. A *compound* parameter consists of one or more terms extending as far as possible to the left or to the right of the macro depending on whether the parameter is “pre-” or “post-”. A *delimited* parameter consists of one or more terms extending as far as possible to the right up to but not including a given token

Table 1: Examples of T_EX and L^AT_EX macros along with their signature.

Macro	Parameters	
	pre	post
overline		[simple]
sqrt		[simple] (T _E X)
		[optional; simple] (L ^A T _E X)
root		[delimited(control(of)); simple]
over, choose	[compound]	[compound]
frac		[simple; simple]
rm, bf, tt, it		[compound]
left		[simple; delimited(control(right)); simple]
sb, sp	[simple]	[simple]
bggroup		[delimited(control(egroup))]
begin		[simple; optional; delimited(control(end)); simple]
proclaim		[token(space); delimited(literal(.)); token(space); delimited(control(par))]

t. An *optional* parameter is either empty or it consists of one or more terms enclosed within a pair of square brackets [and]. The absence of the opening bracket means that the optional parameter is not given. A *token* parameter is a given token *t* representing pure syntactic sugar. It does not properly qualify as a parameter and does not appear in the parsing tree.

Formally tokens and parameter types are defined as follows:

$$\begin{aligned}
 \textit{token} &::= \textit{literal}(v) \mid \textit{space} \mid \textit{control}_{(p_1, p_2)}(v) \\
 \textit{type} &::= \textit{simple} \mid \textit{compound} \mid \textit{delimited}(t) \\
 &\quad \mid \textit{optional} \mid \textit{token}(t)
 \end{aligned}$$

where $t \in \textit{token}$, $v \in \textit{string}$ is an arbitrary string of Unicode characters, $p_1 \in \{\textit{simple}, \textit{compound}\}^*$ and $p_2 \in \textit{type}^*$ are lists of parameter types for the pre- and post-parameters respectively. Note that pre-parameters can be of type *simple* or *compound* only.

The dictionary is a total map

$$\textit{dictionary} : \textit{string} \mapsto \textit{token}$$

such that for each unknown control sequence v we have $\textit{dictionary}(v) = \textit{control}_{(\[], \[]) } (v)$. Table 1 shows part of a possible dictionary for some T_EX and L^AT_EX commands (mostly for mathematics). Note how it is possible to encode the signature for the `\begin` control sequence, although it is not possible to enforce the constraint that the first and the last parameters must have equal value in order for the construct to be balanced.

4 Parsing

We now come to the problem of building the T_EX parsing tree starting the stream of tokens produced by the lexical analyzer. As we have already pointed

out there is no fixed grammar that we can use to generate the parser automatically: authors are free to introduce new macros and hence new ways of structuring the parse tree. Thus we will build the parser “by hand”. More reasons for writing an ad-hoc parser, namely error recovery and incremental-ity, will be discussed later in this section.

The following grammar captures formally the structure of a T_EX parsing tree, which is the outcome of the parser:

$$\begin{aligned}
 \textit{node} &::= \textit{empty} \\
 &\quad \mid \textit{literal}(v) \quad v \in \textit{string} \\
 &\quad \mid \textit{macro}(v, x) \quad v \in \textit{string}, x \in \textit{param}^* \\
 \textit{param} &::= \{a\} \quad a \in \textit{node}^*
 \end{aligned}$$

Note that a parameter is made of a list of nodes and that literals are strings instead of single characters. The `empty` node is used to denote a missing term when one was expected; its role will be clarified later in this section.

The appendix contains the Document Type Definition for the XML representation of T_EX parsing trees. It is simpler than the T_EXML DTD [7] and we are providing it as mere reference.

4.1 Parsing Functions

Table 2 gives the operational semantics of the parser. In this table only, for each $a \in \textit{node}^*$ we define $a! = [\textit{empty}]$ if $a = []$ and $a! = a$ otherwise. There are four parsing functions: \mathcal{T} for terms, \mathcal{A} for pre-parameters, \mathcal{B} for post-parameters, and \mathcal{C} for delimited sequences of terms. Each parsing function is defined by induction on the structure of its arguments. Axioms (rules with no horizontal line) denote base cases, while inference rules define the value

Table 2: Parsing functions for the simplified \TeX markup.

$\forall d \in \text{token}^*$	$\xrightarrow{\mathcal{T}(d)}$	$:$	$\text{node}^* \times \text{token}^* \rightarrow \text{node}^* \times \text{token}^*$
	$\xrightarrow{\mathcal{A}}$	$:$	$\text{node}^* \times \text{type}^* \rightarrow \text{node}^* \times \text{param}^*$
$\forall d \in \text{token}^*$	$\xrightarrow{\mathcal{B}(d)}$	$:$	$\text{type}^* \times \text{token}^* \rightarrow \text{param}^* \times \text{token}^*$
$\forall d \in \text{token}^*, \forall b \in \text{bool}$	$\xrightarrow{\mathcal{C}(d,b)}$	$:$	$\text{node}^* \times \text{token}^* \rightarrow \text{node}^* \times \text{token}^*$

$$\begin{array}{l}
\text{(T.1)} \quad a, [] \xrightarrow{\mathcal{T}(d)} a, [] \quad \text{(T.2)} \quad a, t :: l \xrightarrow{\mathcal{T}(d)} a, t :: l \quad (t \text{ occurs in } d) \\
\text{(T.3)} \quad a, \text{literal}(v) :: l \xrightarrow{\mathcal{T}(d)} a@[\text{literal}(v)], l \quad \text{(T.4)} \quad \frac{a, l \xrightarrow{\mathcal{T}(d)} a', l'}{a, \text{space} :: l \xrightarrow{\mathcal{T}(d)} a', l'} \\
\text{(T.5)} \quad \frac{a, p_1 \xrightarrow{\mathcal{A}} a', x \quad p_2, l \xrightarrow{\mathcal{B}(d)} y, l'}{a, \text{control}_{\langle p_1, p_2 \rangle}(v) :: l \xrightarrow{\mathcal{T}(d)} a'@[\text{macro}(v, x@y)], l'} \\
\text{(A.1)} \quad a, [] \xrightarrow{\mathcal{A}} a, [] \quad \text{(A.2)} \quad \frac{[], p \xrightarrow{\mathcal{A}} a, x}{[], s :: p \xrightarrow{\mathcal{A}} a, x@[\{\{\text{empty}\}\}]} \\
\text{(A.3)} \quad \frac{a, p \xrightarrow{\mathcal{A}} a', x}{a@[n], \text{simple} :: p \xrightarrow{\mathcal{A}} a', x@[\{\{n\}\}]} \quad \text{(A.4)} \quad \frac{[], p \xrightarrow{\mathcal{A}} a', x}{a, \text{compound} :: p \xrightarrow{\mathcal{A}} a', x@[\{\{a\}\}]} \\
\text{(B.1)} \quad [], l \xrightarrow{\mathcal{B}(d)} [], l \quad \text{(B.2)} \quad \frac{p, t' :: l \xrightarrow{\mathcal{B}(d)} x, a}{\text{token}(t) :: p, t' :: l \xrightarrow{\mathcal{B}(d)} x, a} \dagger \quad (t \neq t') \\
\text{(B.3)} \quad \frac{p, [] \xrightarrow{\mathcal{B}(d)} x, l}{\text{token}(t) :: p, [] \xrightarrow{\mathcal{B}(d)} x, l} \dagger \quad \text{(B.4)} \quad \frac{p, l \xrightarrow{\mathcal{B}(d)} x, l'}{\text{token}(t) :: p, t :: l \xrightarrow{\mathcal{B}(d)} x, l'} \\
\text{(B.5)} \quad \frac{[], l \xrightarrow{\mathcal{T}(d)} a, l' \quad p, l' \xrightarrow{\mathcal{B}(d)} x, l''}{\text{simple} :: p, l \xrightarrow{\mathcal{B}(d)} \{a!\} :: x, l''} \quad \text{(B.6)} \quad \frac{[], l \xrightarrow{\mathcal{C}(d, \text{false})} a, l' \quad p, l' \xrightarrow{\mathcal{B}(d)} x, l''}{\text{compound} :: p, l \xrightarrow{\mathcal{B}(d)} \{a!\} :: x, l''} \\
\text{(B.7)} \quad \frac{p, [] \xrightarrow{\mathcal{B}(d)} x, l}{\text{optional} :: p, [] \xrightarrow{\mathcal{B}(d)} \{\}\} :: x, l \quad \text{(B.8)} \quad \frac{[], l \xrightarrow{\mathcal{C}(\text{literal}(\cdot)) :: d, \text{true}} a, l' \quad p, l' \xrightarrow{\mathcal{B}(d)} x, l''}{\text{optional} :: p, \text{literal}(\cdot) :: l \xrightarrow{\mathcal{B}(d)} \{a\} :: x, l''} \\
\text{(B.9)} \quad \frac{p, t :: l \xrightarrow{\mathcal{B}(d)} x, l'}{\text{optional} :: p, t :: l \xrightarrow{\mathcal{B}(d)} \{\}\} :: x, l' \quad (t \neq \text{literal}(\cdot)) \\
\text{(B.10)} \quad \frac{[], l \xrightarrow{\mathcal{C}(t :: d, \text{true})} a, l' \quad p, l' \xrightarrow{\mathcal{B}(d)} x, l''}{\text{delimited}(t) :: p, l \xrightarrow{\mathcal{B}(d)} \{a!\} :: x, l''} \\
\text{(C.1)} \quad a, [] \xrightarrow{\mathcal{C}(d,b)} a, [] \\
\text{(C.2)} \quad a, t :: l \xrightarrow{\mathcal{C}(t :: d, \text{true})} a, l \quad \text{(C.3)} \quad a, t :: l \xrightarrow{\mathcal{C}(d,b)} a, t :: l \quad (t \text{ occurs in } d) \\
\text{(C.4)} \quad \frac{a, t :: l \xrightarrow{\mathcal{T}(d)} a', l' \quad a', l' \xrightarrow{\mathcal{C}(d,b)} a'', l''}{a, t :: l \xrightarrow{\mathcal{C}(d,b)} a'', l''} \quad (t \notin d)
\end{array}$$

of a parsing function (the conclusion, below the line) in terms of the value of one or more recursive calls to other functions (the premises, above the line). Right arrows denote the action of parsing. Arrows are decorated with a label that identifies the parser along with its parameters, if any. The \mathcal{T} , \mathcal{B} , and \mathcal{C} parsers have a parameter representing the list of delimiters in the order they are expected, with the head of the list being the first expected delimiter. The \mathcal{C} parser also has a boolean parameter indicating whether the parser should or should not “eat” the delimiter when it is eventually met.

The root parsing function is \mathcal{T} . Given a delimiter $t \in \textit{token}$ and a token stream $l \in \textit{token}^*$ we have

$$\square, l \xrightarrow{\mathcal{T}([t])} [n], l'$$

where $n \in \textit{node}$ is the parsed term and $l' \in \textit{token}^*$ is the part of the token stream that has not been consumed. Spaces are ignored when parsing terms and pre-parameters (rule *T.4*), but not when parsing post-parameters (rule *B.4*). The \mathcal{A} function differs from the other parsing functions because by the time a macro with pre-parameters is encountered, pre-parameters have already been parsed. The lists $a \in \textit{node}^*$ in the \mathcal{T} , \mathcal{A} , and \mathcal{C} parsers represent the terms accumulated before the term being parsed. Note that pre-parameters are inserted at the end of the parameter list (rules *A.2* to *A.4*) and that post-parameters are inserted at the beginning of the parameter list (rules *B.5* to *B.10*). This way parameter nodes appear in the parse tree in the same order as in the original token stream (rule *T.5*).

4.1.1 Example

Given that the input buffer contains the T_EX source shown in Fig. 1, the lexical analyzer would produce the following stream of tokens:

$$l_0 \stackrel{\text{def}}{=} [\text{control}_{\langle \square, [\text{delimited}(\text{control}(\text{egroup})) \rangle]}(\text{bgroup}); \\ \text{literal}(1); \text{control}_{\langle [\text{compound}], [\text{compound}] \rangle}(\text{over}); \\ \text{control}_{\langle \square, [\text{delimited}(\text{control}(\text{egroup})) \rangle]}(\text{bgroup}); \\ \text{literal}(x); \text{literal}(+); \text{literal}(1); \\ \text{control}(\text{egroup}); \text{control}_{\langle [\text{simple}], [\text{simple}] \rangle}(\text{sp}); \\ \text{literal}(2); \text{control}(\text{egroup})]$$

By the application of the parsing rules given in Table 2 it can be shown that

$$\square, l_0 @ [\text{control}(\text{eoi})] \xrightarrow{\mathcal{T}([\text{control}(\text{eoi})])} [n], [\text{control}(\text{eoi})]$$

where $n \in \textit{node}$ is the same tree shown in Fig. 1 except that the \mathbf{g} nodes are labeled with \mathbf{bgroup} .

4.2 Error Recovery

Parsing functions are all total functions, they always produce a result, even when the input token

stream is malformed. Unlike parsers of batch T_EX converters or the T_EX parser itself, there will often be moments during the editing process when the input buffer contains incorrect or incomplete markup, for example because not all the required parameters of a macro have been entered yet. The parser must recover from such situations in a tolerant and hopefully sensible way. We distinguish three kinds of situations: *missing parameters*, *pattern mismatch*, and *ambiguity*, which we examine in the rest of this section.

4.2.1 Missing Parameters

Consider an input token stream representing the sole $\backslash\text{over}$ macro with no arguments provided:

$$l_1 \stackrel{\text{def}}{=} [\text{control}_{\langle [\text{compound}], [\text{compound}] \rangle}(\text{over}); \\ \text{control}(\text{eoi})]$$

It is easy to check that

$$\square, l_1 \xrightarrow{\mathcal{T}([\text{control}(\text{eoi})])} [\text{macro}(\text{over}, [\text{empty}; \text{empty}]), \\ \text{control}(\text{eoi})]$$

More generally the parser inserts **empty** nodes in the parsing tree wherever an expected parameter is not found in the token stream. This behavior can be seen in rule *A.2* and also in rules *B.5*, *B.6*, and *B.10* where the $!$ operator is used. For optional parameters an empty node list is admitted (rules *B.7* and *B.8*).

The presence of **empty** nodes guarantees that the generated tree is structurally well-formed, which is crucial for the subsequent transformation phase. It also allows the application to give the user feedback indicating the absence of required parameters. In the example above, for instance, the application may display something like \square suggesting that a fraction was entered, but neither the numerator nor the denominator have been.

4.2.2 Pattern Mismatch

Rules *B.2* and *B.3* have been marked with a \dagger to indicate that the parser expects a token which is not found in the token stream. In both cases the parser will typically notify the user with a warning message.

4.2.3 Ambiguities

In T_EX one cannot pass a macro with parameters as the parameter of another macro, unless the parameter is enclosed within a group. For example, it is an error to write $\backslash\text{sqrt}\backslash\text{sqrt}\{x\}$, the correct form is $\backslash\text{sqrt}\{\backslash\text{sqrt}\{x\}\}$. Because we treat the left curly brace like any other macro, grouping would not help our parser in resolving ambiguities. However, the

parser knows how many parameters a macro needs, because the token representing the control sequence has been annotated with such information by the lexer. When processing a macro with arguments the parser behaves “recursively”, it does not let an incomplete macro to be “captured” if it was not passed as parameter of an outer macro. A consequence of this extension is that any well-formed fragment of \TeX markup is accepted by our parser resulting in the same structure, but there are some strings accepted by our parser that cause the \TeX parser to fail.

4.3 Incremental Parsing

Parsing must be efficient because it is performed in real-time, in principle at every modification of the input buffer, no matter how simple the modification is. Fortunately \TeX markup exhibits good *locality*, that is small modifications in the document cause small modifications in the parsing tree. Consequently we can avoid re-parsing the whole source document, we just need to re-parse a small interval of the input buffer around the point where the modification has occurred, and adjust the parsing tree accordingly. Let us consider again the example of Fig. 1 and suppose that a change is made in the markup

$$\{1\over{1+x}^2\} \Rightarrow \{1\over{1+x+y}^2\}$$

(a $+y$ is added to the denominator of the fraction). To be conservative we can re-parse the smallest term within braces that includes the modified part (the underlined fragments). Once the term has been re-parsed it has to be substituted in place of the old term in the parsing tree.

In order to compute the interval of the input buffer to be re-parsed we annotate the nodes of the parsing tree with information about the first and the last characters of the buffer which were scanned while building the node and all of its children. A simple visit of the tree can locate the smaller interval affected by the modification.

Curly braces occur frequently enough in the markup to give good granularity for re-parsing. At the same time limiting re-parsing to braced terms helps control the costs related to the visit to the parsing tree and to the implementation of the incremental parsing and transformation machinery.

5 Transformation

The transformation phase recognizes structured patterns in the parsing tree and generates corresponding fragments of the result document. We have already anticipated that XSLT is a very natural choice

for the implementation of this phase. Besides, XSLT stylesheets can be extended very easily, by providing new *templates* that recognize and properly handle new macros that an author has introduced.

We can see in Fig. 2 two sample templates taken from an XSLT stylesheet for converting the internal parsing tree into a MathML tree. Both templates have a preamble made of an `xs1:if` construct which we will discuss later in this section. Since the \TeX tree and the MathML tree are almost isomorphic (Fig. 1) the transformation is generally very simple and in many cases it amounts at just renaming the node labels. Template (a) is one such case: it matches any node in the parsing tree with label `macro` and having the `name` attribute set to `over`. The node for the `\over` macro corresponds naturally to the `mfrac` element in MathML. The two parameters of `\over` are transformed recursively by applying the stylesheet templates to the first and second child nodes (`p[1]` means “the first `p` child of this node”, similarly `p[2]` refers to the second `p` child).

Template (b) is slightly more complicated and shows one case where there is some change in the structure. For combined sub/super scripts \TeX accepts a sequence of `_` and `^` no matter in what order they occur, but MathML has a specific element for such expressions, namely `msubsup`. The template matches an `sb` node whose first parameter contains an `sp` node, thus detecting a `...^..._...` fragment of markup, then the corresponding `msubsup` element is created and its three children accessed in the proper position of the parsing tree. A symmetric template will handle the case where the subscript occurs before the superscript.

5.1 Incremental Transformation

As we have done for parsing, for transformations we also need to account for their cost. In a batch, one-shot conversion from \TeX this is not generally an issue, but in an interactive authoring tool a transformation is required at every modification of the parsing tree in order to update the view of the document.

Intuitively, we can reason that if only a fragment of the parsing tree has changed, we need re-transform only that fragment and substitute the result in the final document. This technique makes two assumptions: (1) that transformations are context-free; that is, the transformation of a fragment in the parsing tree is not affected by the context in which the fragment occurs; (2) that we are able to relate corresponding fragments between the parsing and the result trees.

<pre> <xsl:template match="macro[@name='over']"> <m:frac> <xsl:if test="@id"> <xsl:attribute name="xref"> <xsl:value-of select="@id"/> </xsl:attribute> </xsl:if> <xsl:apply-templates select="p[1]"/> <xsl:apply-templates select="p[2]"/> </m:frac> </xsl:template> </pre>	<pre> <xsl:template match="macro[@name='sb'] [p[1]/*[1][self::macro[@name='sp']]]"> <m:msubsup> <xsl:if test="@id"> <xsl:attribute name="xref"> <xsl:value-of select="@id"/> </xsl:attribute> </xsl:if> <xsl:apply-templates select="p[1]/*p[1]"/> <xsl:apply-templates select="p[2]"/> <xsl:apply-templates select="p[1]/*p[2]"/> </m:msubsup> </xsl:template> </pre>
(a)	(b)

Figure 2: Example of XSLT templates for the transformation of the internal parsing tree into a MathML tree. MathML elements can be distinguished because of the `m:` prefix.

Template (b) in Fig. 2 shows one case where the transformation is not context free: the deeper `sp` node is not processed as if it would occur alone, but it is “merged” together with its parent. More generally we can imagine that transformations can make almost arbitrary re-arrangements of the structure. This problem cannot be solved unless we make some assumptions, and the one we have already committed to in Sect. 4 is that braces define “black-box” fragments which can be transformed in isolation, without context dependencies.

As for the matter of relating corresponding fragments of the two documents, we use identifiers and references. Each node in the parsing tree is annotated with a unique identifier (in our sample templates we are assuming that the identifier is a string in the `id` attribute). Templates create corresponding `xref` attributes in the result document “pointing” to the fragment with the same identifier in the parsing tree. This way, whenever a fragment of the parsing tree is re-transformed, it replaces the fragment in the result document with the same identifier.

More generally, back-pointers provide a mechanism for relating the view of the document with the source markup. This way it is possible to perform operations like selection or cut-and-paste that, while having a visual effect in the view, act indirectly at the content/markup level.

6 Conclusion

We have presented architectural and implementation issues of an interactive editor based on \TeX syntax which allows flexible customization and content-oriented authoring. \TeX macs² is probably the existing application that most closely adopts such architecture, with the difference that \TeX macs does not

stick to \TeX syntax as closely as we do and that, apart from being a complete (and cumbersome) editing tool and not just an interface, it uses encoding and transformation technologies not based on standard languages (XML [2] and XSLT [3]).

Among batch conversion tools we observe a tendency to move towards the processing of content. The \TeX to MathML converter by Igor Rodionov and Stephen Watt at the University of Western Ontario [8, 9] is one such tool, and the recent Hermes converter by Romeo Anghelache [10] is another. These represent significant steps forwards when compared to converters such as \LaTeX 2HTML.³

A prototype tool called Edi \TeX , based on the architecture described in this paper, has been developed and is freely available along with its source code.⁴ No mention of MathML is made in the name of the tool to remark the fact that the architecture is very general and can be adapted to other kinds of markup. The prototype is currently being used as interface for a proof-assistant application where editing of complex mathematical formulas and proofs is required. In this respect we should remark that \TeX syntax is natural for “real” mathematics, but it quickly becomes clumsy when used for writing terms of programming languages or λ -calculus. This is mainly due to the conventions regarding spaces (for instance, spaces in the λ -calculus denote function application) and identifiers (the rule “one character is one identifier” is fine for mathematics, but not for many other languages). Note however that, since the lexical analyzer is completely separate from the rest of the architecture, the token stream being its interface, it can be easily targeted to a language with different conventions than those of \TeX .

² <http://www.texmacs.org/>

³ <http://www.latex2html.org/>

⁴ <http://helm.cs.unibo.it/software/editex/>

The idea of using some sort of restricted \TeX syntax for representing mathematical expressions is not new. For example, John Forkosh's \MimeTeX ⁵ generates bitmap images of expressions to be embedded in Web pages. However, to the best of our knowledge the formal specification of the parser for simplified \TeX markup presented in Sect. 4 is unique of its kind. A straightforward implementation based directly on the rules given in Table 2 amounts at only just 70 lines of functional code (in an ML dialect), which can be considered something of an achievement given that parsing \TeX is normally regarded as a hard task. By comparison, the parsing code in \MimeTeX amounts to nearly 350 lines of C code after stripping away the comments.

One may argue that the simplified \TeX markup is too restrictive, but in our view this is just the sensible fragment of \TeX syntax that the average user should be concerned about. In fact the remaining syntactic expressiveness provided by \TeX is mainly required for the implementation of complex macros and of system internals, which should never surface at the document level. By separating the transformation phase we shift the mechanics of macro expansion to a different level which can be approached with different (more appropriate) languages. Since this mode of operation makes the system more flexible we believe that our design is a valuable contribution which may provide an architecture for other implementers to adopt.

References

- [1] The Unicode Consortium: The Unicode Standard, Version 4.0, Boston, MA, Addison-Wesley (2003). <http://www.unicode.org/>
- [2] Tim Bray, Jean Paoli, C.M. Sperberg-McQueen, Eve Maler (editors): Extensible Markup Language (XML) 1.0 (2nd Edition), W3C Recommendation (2000). <http://www.w3.org/TR/2000/REC-xml-20001006>
- [3] James Clark (editor): XML Transformations (XSLT) Version 1.0, W3C Recommendation (1999). <http://www.w3.org/TR/1999/REC-xslt-19991116>
- [4] Ron Ausbrooks, Stephen Buswell, Stéphane Dalmas, Stan Devitt, Angel Diaz, et al.: Mathematical Markup Language (MathML) Version 2.0 (2nd Edition) W3C Recommendation, (2003). <http://www.w3.org/TR/2003/REC-MathML2-20031021/>
- [5] Donald E. Knuth: The \TeX book, Addison-Wesley, Reading, MA, USA (1994).
- [6] Leslie Lamport: A Document Preparation System: \LaTeX , Addison-Wesley, Reading, MA, USA (1986).
- [7] Douglas Lovell: \TeX XML: Typesetting XML with \TeX , *TUGboat*, 20(3), pp. 176–183 (September 1999).
- [8] Sandy Huerter, Igor Rodionov, Stephen M. Watt: Content-Faithful Transformations for MathML, Proc. International Conference on MathML and Math on the Web (MathML 2002), Chicago, USA (2002). <http://www.mathmlconference.org/2002/presentations/huerter/>
- [9] Stephen M. Watt: Conserving implicit mathematical semantics in conversion between \TeX and MathML, *TUGboat*, 23(1), pp. 108–108 (2002).
- [10] Romeo Anghelache: \LaTeX -based authoring tool, Deliverable D4.d, MoWGLI Project (2003). <http://relativity.livingreviews.org/Info/AboutLR/mowgli/index.html>

Appendix: The TML DTD

```

<!ENTITY % TML.node "
  empty|space|literal|macro">
<!ENTITY % TML.common.attrib "
  id          CDATA #IMPLIED
  xref        CDATA #IMPLIED
  start       NMTOKEN #IMPLIED
  end         NMTOKEN #IMPLIED">
<!ELEMENT empty EMPTY>
<!ATTLIST empty %TML.common.attrib;>
<!ELEMENT space EMPTY>
<!ATTLIST space
  %TML.common.attrib
  name      NMTOKEN #IMPLIED
  literal   CDATA #IMPLIED>
<!ELEMENT literal #PCDATA>
<!ATTLIST literal
  %TML.common.attrib;
  name      NMTOKEN #IMPLIED>
<!ELEMENT macro (p)*>
<!ATTLIST macro
  %TML.common.attrib;
  name      NMTOKEN #REQUIRED
  literal   CDATA #IMPLIED>
<!ELEMENT p (%TML.node;)*>
<!ATTLIST p %TML.common.attrib;>

```

⁵ <http://www.ctan.org/tex-archive/support/mimetex/>