

Semantic Markup in T_EX/L^AT_EX*

Michael Kohlhase
Jacobs University, Bremen
<http://kwarc.info/kohlhase>

April 3, 2016

Abstract

We present a collection of T_EX macro packages that allow to markup T_EX/L^AT_EX documents semantically without leaving the document format, essentially turning T_EX/L^AT_EX into a document format for mathematical knowledge management (MKM).

Contents

1	Introduction	2
1.1	The XML vs. T _E X/L ^A T _E X Formats and Workflows	2
1.2	A L ^A T _E X-based Workflow for XML-based Mathematical Documents	4
2	The Packages of the sT_EX Collection	4
2.1	Content Markup of Mathematical Formulae in T _E X/L ^A T _E X	4
2.2	Mathematical Statements	7
2.3	Context Markup for Mathematics	7
2.4	Mathematical Document Classes	7
2.5	Conclusion	7
2.6	Licensing, Download and Setup	8
3	Utilities	8
4	The Implementation	9
4.1	Package Options	9
4.2	The sT _E X Logo	9

*Version ? (last revised ?)

1 Introduction

The last few years have seen the emergence of various content-oriented XML-based, content-oriented markup languages for mathematics on the web, e.g. OpenMath [Bus+04], content MathML [Aus+03], or our own OMDoc [Koh06]. These representation languages for mathematics, that make the structure of the mathematical knowledge in a document explicit enough that machines can operate on it. Other examples of content-oriented formats for mathematics include the various logic-based languages found in automated reasoning tools (see [RV01] for an overview), program specification languages (see e.g. [Ber89]).

The promise if these content-oriented approaches is that various tasks involved in “doing mathematics” (e.g. search, navigation, cross-referencing, quality control, user-adaptive presentation, proving, simulation) can be machine-supported, and thus the working mathematician is relieved to do what humans can still do infinitely better than machines: The creative part of mathematics — inventing interesting mathematical objects, conjecturing about their properties and coming up with creative ideas for proving these conjectures. However, before these promises can be delivered upon (there is even a conference series [] studying “Mathematical Knowledge Management (MKM)”), large bodies of mathematical knowledge have to be converted into content form.

Even though MathML is viewed by most as the coming standard for representing mathematics on the web and in scientific publications, it has not not fully taken off in practice. One of the reasons for that may be that the technical communities that need high-quality methods for publishing mathematics already have an established method which yields excellent results: the \TeX / \LaTeX system: and a large part of mathematical knowledge is prepared in the form of \TeX / \LaTeX documents.

\TeX [Knu84] is a document presentation format that combines complex page-description primitives with a powerful macro-expansion facility, which is utilized in \LaTeX (essentially a set of \TeX macro packages, see [Lam94]) to achieve more content-oriented markup that can be adapted to particular tastes via specialized document styles. It is safe to say that \LaTeX largely restricts content markup to the document structure¹, and graphics, leaving the user with the presentational \TeX primitives for mathematical formulae. Therefore, even though \LaTeX goes a great step into the direction of an MKM format, it is not, as it lacks infrastructure for marking up the functional structure of formulae and mathematical statements, and their dependence on and contribution to the mathematical context.

1.1 The XML vs. \TeX / \LaTeX Formats and Workflows

MathML is an XML-based markup format for mathematical formulae, it is standardized by the World Wide Web Consortium in [Aus+03], and is supported by the major browsers. The MathML format comes in two integrated components: presentation MathML presentation MathML and content MathML content MathML.

¹supplying macros e.g. for sections, paragraphs, theorems, definitions, etc.

The former provides a comprehensive set of layout primitives for presenting the visual appearance of mathematical formulae, and the second one the functional/-logical structure of the conveyed mathematical objects. For all practical concerns, presentation MathML is equivalent to the math mode of \TeX . The text mode facilities of \TeX (and the multitude of \LaTeX classes) are relegated to other XML formats, which embed MathML.

The programming language constructs of \TeX (i.e. the macro definition facilities²) are relegated to the XML programming languages that can be used to develop language extensions. transformation language XSLT [Dea99; Kay00] or proper XML-enabled The XML-based syntax and the separation of the presentational-, functional- and programming/extensibility concerns in MathML has some distinct advantages over the integrated approach in $\text{\TeX}/\text{\LaTeX}$ on the services side: MathML gives us better

- integration with web-based publishing,
- accessibility to disabled persons, e.g. (well-written) MathML contains enough structural information to supports screen readers.
- reusability, searchability and integration with mathematical software systems (e.g. copy-and-paste to computer algebra systems), and
- validation and plausibility checking.

On the other hand, $\text{\TeX}/\text{\LaTeX}$'s adaptable syntax and tightly integrated programming features within has distinct advantages on the authoring side:

- The $\text{\TeX}/\text{\LaTeX}$ syntax is much more compact than MathML (see the difference in Figure 1 and Equation 1), and if needed, the community develops \LaTeX packages that supply new functionality in with a succinct and intuitive syntax.
- The user can define ad-hoc abbreviations and bind them to new control sequences to structure the source code.
- The $\text{\TeX}/\text{\LaTeX}$ community has a vast collection of language extensions and best practice examples for every conceivable publication purpose and an established and very active developer community that supports these.
- There is a host of software systems centered around the $\text{\TeX}/\text{\LaTeX}$ language that make authoring content easier: many editors have special modes for \LaTeX , there are spelling/style/grammar checkers, transformers to other markup formats, etc.

In other words, the technical community is heavily invested in the whole workflow, and technical know-how about the format permeates the community.

²We count the parser manipulation facilities of \TeX , e.g. category code changes into the programming facilities as well, these are of course impossible for MathML, since it is bound to XML syntax.

Since all of this would need to be re-established for a MathML-based workflow, the technical community is slow to take up MathML over $\text{\TeX}/\text{\LaTeX}$, even in light of the advantages detailed above.

1.2 A \LaTeX -based Workflow for XML-based Mathematical Documents

An elegant way of sidestepping most of the problems inherent in transitioning from a \LaTeX -based to an XML-based workflow is to combine both and take advantage of the respective advantages.

The key ingredient in this approach is a system that can transform $\text{\TeX}/\text{\LaTeX}$ documents to their corresponding XML-based counterparts. That way, XML-documents can be authored and prototyped in the \LaTeX workflow, and transformed to XML for publication and added-value services, combining the two workflows.

There are various attempts to solve the $\text{\TeX}/\text{\LaTeX}$ to XML transformation problem (see [Sta+09] for an overview); the most mature is probably Bruce Miller’s $\text{\LaTeX}XML$ system [LTX]. It consists of two parts: a re-implementation of the \TeX analyzer with all of its intricacies, and an extensible XML emitter (the component that assembles the output of the parser). Since the \LaTeX style files are (ultimately) programmed in \TeX , the \TeX analyzer can handle all \TeX extensions, including all of \LaTeX . Thus the $\text{\LaTeX}XML$ parser can handle all of $\text{\TeX}/\text{\LaTeX}$, if the emitter is extensible, which is guaranteed by the $\text{\LaTeX}XML$ binding language: To transform a $\text{\TeX}/\text{\LaTeX}$ document to a given XML format, all \TeX extensions³ must have “ $\text{\LaTeX}XML$ bindings” binding, i.e. a directive to the $\text{\LaTeX}XML$ emitter that specifies the target representation in XML.

2 The Packages of the \sTeX Collection

In the following, we will shortly preview the packages and classes in the \sTeX collection. They all provide part of the solution of representing semantic structure in the $\text{\TeX}/\text{\LaTeX}$ workflow. We will group them by the conceptual level they address¹

2.1 Content Markup of Mathematical Formulae in $\text{\TeX}/\text{\LaTeX}$

The first two packages are concerned basically with the math mode in \TeX , i.e. mathematical formulae. The underlying problem is that run-of-the-mill $\text{\TeX}/\text{\LaTeX}$ only specifies the presentation (i.e. what formulae look like) and not their content (their functional structure). Unfortunately, there are no good methods (yet) to infer the latter from the former, but there are ways to get presentation from content.

³i.e. all macros, environments, and syntax extensions used in the source document

¹EdNOTE: come up with a nice overview figure here!

Consider for instance the following “standard notations”⁴ for binomial coefficients: $\binom{n}{k}$, ${}_nC^k$, C_k^n all mean the same thing: $\frac{n!}{k!(n-k)!}$. This shows that we cannot hope to reliably recover the functional structure (in our case the fact that the expression is constructed by applying the binomial function to the arguments n and k) from the presentation alone.

The solution to this problem is to dump the extra work on the author (after all she knows what she is talking about) and give them the chance to specify the intended structure. The markup infrastructure supplied by the \S T E X collection lets the author do this without changing⁵ the visual appearance, so that the \La T E X workflow is not disrupted. . We speak of semantic preloading for this process and call our collection of macro packages \S T E X (Semantic \T E X). For instance, we can now write

$\backslash\text{CSumLimits}\{k\}1\backslash\infty\{\backslash\text{Cexp}\{x\}k\}$ instead of the usual $\backslash\text{sum}_{k=1}^{\infty} x^k$ (1)

In the first form, we specify that you are applying a function ($\text{CSumLimits} \hat{=}$ Sum with Limits) to 4 arguments: *i*) the bound variable k (that runs from) *ii*) the number 1 (to) *iii*) ∞ (to infinity summing the terms) *iv*) $\backslash\text{Cexp}\{x\}k$ (i.e. x to the power k). In the second form, we merely specify that \La T E X should draw a capital Sigma character (Σ) with a lower index which is an equation $k = 1$ and an upper index ∞ . Then it should place next to it an x with an upper index k .

Of course human readers (that understand the math) can infer the content structure from this presentation, but the \La T E X ML converter (who does not understand the math) cannot, but we want to have the content MathML expression in Figure 1.

```
% <math xmlns="http://www.w3.org/1998/Math/MathML">
%   <bind>
%     <apply><sumlimits/><cn>1</cn><infinity/></apply>
%     <bvar><ci>k</ci></bvar>
%     <apply><exp/><ci>x</ci><ci>k</ci></apply>
%   </bind>
% </math>
%
```

Example 1: Content MathML Form of $\sum_{k=1}^{\infty} x^k$

Obviously, a converter can infer this from the first \La T E X structure with the help of the curly braces that indicate the argument structure, but not from the second (because it does not understand the math).

The nice thing about the cmathml infrastructure is that you can still run \La T E X over the first form and get the same formula in the DVI file that you would have gotten from running it over the second form. That means, if the author is prepared

⁴The first one is standard e.g. in Germany and the US, and the last one in France

⁵However, semantic annotation will make the author more aware of the functional structure of the document and thus may in fact entice the author to use presentation in a more consistent way than she would usually have.

to write the mathematical formulae a little differently in her \LaTeX sources, then she can use them in XML and \LaTeX at the same time.

2.1.1 cmathml: Encoding Content MathML in \TeX / \LaTeX

The `cmathml` package (see [Koh15a]) provides a set of macros that correspond to the K-14 fragment of mathematics (Kindergarten to undergraduate college level ($\hat{=}$ 14th grade)). We have already seen an example above in equation (1), where the content markup in \TeX corresponds to a content MathML-expression (and can actually be translated to this by the \LaTeX XML system.) However, the content MathML vocabulary is fixed in the MathML specification and limited to the K-14 fragment; the notation of mathematics of course is much larger and extensible on the fly.

2.1.2 presentation: Flexible Presentation for Semantic Macros

The `presentation` package (see [KG15]) supplies an infrastructure that allows to specify the presentation of semantic macros, including preference-based bracket elision. This allows to markup the functional structure of mathematical formulae without having to lose high-quality human-oriented presentation in \LaTeX . Moreover, the notation definitions can be used by MKM systems for added-value services, either directly from the \LaTeX sources, or after translation.

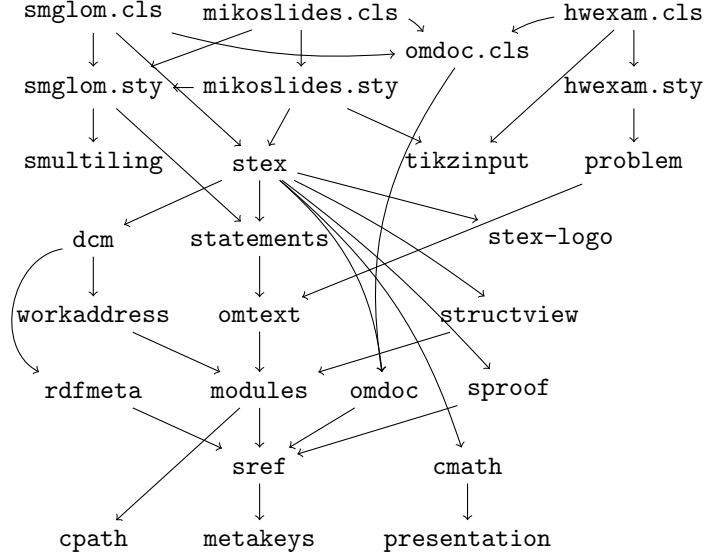


Figure 1: The \LaTeX packages and their dependencies.

2.2 Mathematical Statements

2.2.1 `statements`: Extending Content Macros for Mathematical Notation

The `statements` package (see[Koh15c]) provides semantic markup facilities for mathematical statements like Theorems, Lemmata, Axioms, Definitions, etc. in \LaTeX files. This structure can be used by MKM systems for added-value services, either directly from the \LaTeX sources, or after translation.

2.2.2 `sproof`: Extending Content Macros for Mathematical Notation

The `sproof` package (see [Koh15b]) supplies macros and environment that allow to annotate the structure of mathematical proofs in \LaTeX files. This structure can be used by MKM systems for added-value services, either directly from the \LaTeX sources, or after translation.

2.3 Context Markup for Mathematics

2.3.1 `modules`: Extending Content Macros for Mathematical Notation

The `modules` package (see [KGA15]) supplies a definition mechanism for semantic macros and a non-standard scoping construct for them, which is oriented at the semantic dependency relation rather than the document structure. This structure can be used by MKM systems for added-value services, either directly from the \LaTeX sources, or after translation.

2.4 Mathematical Document Classes

2.4.1 `OMDoc` Documents

The `omdoc` package provides an infrastructure that allows to markup OMDoc documents in \LaTeX . It provides `omdoc.cls`, a class with the and `omdocdoc.sty`

2.5 Conclusion

The \LaTeX collection provides a set of semantic macros that extends the familiar and time-tried \LaTeX workflow in academics until the last step of Internet publication of the material. For instance, an SMGloM module can be authored and maintained in \LaTeX using a simple text editor, a process most academics in technical subjects are well familiar with. Only in a last publishing step (which is fully automatic) does it get transformed into the XML world, which is unfamiliar to most academics.

Thus, \LaTeX can serve as a conceptual interface between the document author and MKM systems: Technically, the semantically preloaded \LaTeX documents are transformed into the (usually XML-based) MKM representation formats, but conceptually, the ability to semantically annotate the source document is sufficient.

The \LaTeX macro packages have been validated together with a case study [Koh04], where we semantically preload the course materials for a two-semester course in

Computer Science at Jacobs University Bremen and transform them to the OM-
Doc MKM format.

2.6 Licensing, Download and Setup

The $\text{\texttt{sTeX}}$ packages are licensed under the $\text{\texttt{L\^AT\^EX}}$ Project Public License [Pro07], which basically means that they can be downloaded, used, copied, and even modified by anyone under a set of simple conditions (e.g. if you modify you have to distribute under a different name).

The $\text{\texttt{sTeX}}$ packages and classes are available from the Comprehensive $\text{\texttt{T\^E\textsubscript{X}}}$ Archive Network (CTAN [CTAN]) and are part of the primary $\text{\texttt{T\^E\textsubscript{X}}}$ / $\text{\texttt{L\^AT\^EX}}$ distributions (e.g. $\text{\texttt{TeXlive}}$ [] and $\text{\texttt{MikTeX}}$ []). The development version is on GitHub [sTeX], it can be cloned or forked from the repository URL

`https://github.com/KWARC/sTeX.git`

It is usually a good idea to enlarge the internal memory allocation of the $\text{\texttt{T\^E\textsubscript{X}}}$ / $\text{\texttt{L\^AT\^EX}}$ executables. This can be done by adding the following configurations in `texmf.cnf` (or changing them, if they already exist). Note that you will probably need `sudo` to do this.

```
% max_in_open = 50           % simultaneous input files and error insertions ,
% param_size = 20000         % simultaneous macro parameters , also applies to MF
% nest_size = 1000          % simultaneous semantic levels (e.g., groups)
% stack_size = 10000         % simultaneous input sources
% main_memory = 12000000
%
```

After that, you have to run the

```
% sudo fmtutil-sys --all
%
```

3 Utilities

To simplify dealing with $\text{\texttt{sTeX}}$ documents, we are providing a small collection of command line utilities, which we will describe here. For details and downloads go to <http://kwarc.info/projects/stex>.

`msplit` splits an $\text{\texttt{sTeX}}$ file into smaller ones (one module per file)

`rf` computes the “reuse factor”, i.e. how often $\text{\texttt{sTeX}}$ modules are reused over a collection of documents

`sgraph` visualizes the module graph

`sms` computes the $\text{\texttt{sTeX}}$ module signatures for a give $\text{\texttt{sTeX}}$ file

`bms` proposes a sensible module structure for an un-annotated $\text{\texttt{sTeX}}$ file

4 The Implementation

4.1 Package Options

The first step is to declare (a few) package options that handle whether certain information is printed or not. They all come with their own conditionals that are set by the options.

```
1 <*package>
2 \DeclareOption*{\PassOptionsToPackage{\CurrentOption}{statements}}
3                                     \PassOptionsToPackage{\CurrentOption}{structview}
4                                     \PassOptionsToPackage{\CurrentOption}{sproofs}
5                                     \PassOptionsToPackage{\CurrentOption}{omdoc}
6                                     \PassOptionsToPackage{\CurrentOption}{cmath}
7                                     \PassOptionsToPackage{\CurrentOption}{dcm}}
8 \ProcessOptions
```

Then we make sure that the necessary packages are loaded (in the right versions).

```
9 \RequirePackage{stex-logo}
10 \RequirePackage{statements}
11 \RequirePackage{structview}
12 \RequirePackage{sproof}
13 \RequirePackage{omdoc}
14 \RequirePackage{cmath}
15 \RequirePackage{dcm}
16 </package>
```

4.2 The \TeX Logo

To provide default identifiers, we tag all elements that allow `xml:id` attributes by executing the `numberIt` procedure from `omdoc.sty.ltxml`.

```
17 <*logo>
18 \RequirePackage{xspace}
19 \def\stex{\@ifundefined{texorpdfstring}{\let\texorpdfstring\@firstoftwo}{\texorpdfstring{\rais
20 \def\sTeX{\stex}
21 </logo>
```