

Semantic Markup in T_EX/L^AT_EX

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

November 23, 2017

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	4
1.1	The XML vs. T _E X/L ^A T _E X Formats and Workflows	4
1.2	A L ^A T _E X-based Workflow for XML-based Mathematical Documents	6
1.3	Generating OMDoc from sT _E X	6
1.4	Conclusion	6
1.5	Licensing, Download and Setup	7
2	The Packages of the sT_EX Collection	9
2.1	The sT _E X Distribution	9
2.2	Content Markup of Mathematical Formulae in T _E X/L ^A T _E X	9
2.3	Mathematical Statements	10
2.4	Context Markup for Mathematics	10
2.5	Mathematical Document Classes	11
2.6	Metadata	11
2.7	Support for MathHub	12
2.8	Auxiliary Packages	12
3	Workflows and Best Practices	13
3.1	The “Little Modules” Approach	13
3.2	Basic Utilities & Makefiles	14
3.3	MathHub: a Portal for Active Mathematical Documents	14
3.4	lmh: MathHub’s Build System Locally	14

4	The Implementation	16
4.1	Package Options	16
4.2	The $\text{\texttt{STEX}}$ Logo	16

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 [BusCapCar:2oms04], content MathML [CarIon:MathML03], or our own OMDoc [Kohlhase:OMDoc1.2]. 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 [RobVor:hoar01] for an overview), program specification languages (see e.g. [Bergstra:as89]).

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 [MKM-IG-Meetings:online] 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 [Knuth:ttb84] 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 [Lamport:ladps94]) 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 stan-

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

dardized by the World Wide Web Consortium in [CarIon:MathML03], and is supported by the major browsers. The MathML format comes in two integrated components: presentation MathML presentation MathML and content MathML content MathML. 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 $\text{T}_{\text{E}}\text{X}$. The text mode facilities of $\text{T}_{\text{E}}\text{X}$ (and the multitude of $\text{L}^{\text{A}}\text{T}_{\text{E}}\text{X}$ classes) are relegated to other XML formats, which embed MathML.

The programming language constructs of $\text{T}_{\text{E}}\text{X}$ (i.e. the macro definition facilities²) are relegated to the XML programming languages that can be used to develop language extensions. transformation language XSLT [Deach:exls99; Kay:xpr00] 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 $\text{T}_{\text{E}}\text{X}/\text{L}^{\text{A}}\text{T}_{\text{E}}\text{X}$ 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, searchabiliby and integration with mathematical software systems (e.g. copy-and-paste to computer algebra systems), and
- validation and plausibility checking.

On the other hand, $\text{T}_{\text{E}}\text{X}/\text{L}^{\text{A}}\text{T}_{\text{E}}\text{X}$ /s adaptable syntax and tightly integrated programming features within has distinct advantages on the authoring side:

- The $\text{T}_{\text{E}}\text{X}/\text{L}^{\text{A}}\text{T}_{\text{E}}\text{X}$ syntax is much more compact than MathML, and if needed, the community develops $\text{L}^{\text{A}}\text{T}_{\text{E}}\text{X}$ 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 $\text{T}_{\text{E}}\text{X}/\text{L}^{\text{A}}\text{T}_{\text{E}}\text{X}$ 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 $\text{T}_{\text{E}}\text{X}/\text{L}^{\text{A}}\text{T}_{\text{E}}\text{X}$ language that make authoring content easier: many editors have special modes for $\text{L}^{\text{A}}\text{T}_{\text{E}}\text{X}$, there are spelling/style/grammar checkers, transformers to other markup formats, etc.

²We count the parser manipulation facilities of $\text{T}_{\text{E}}\text{X}$, 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.

In other words, the technical community is heavily invested in the whole workflow, and technical know-how about the format permeates the community. 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 [StaGinDav:maacl09] for an overview); the most mature is probably Bruce Miller’s \LaTeX XML system [Miller:latexml:online]. 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 \LaTeX XML parser can handle all of $\text{\TeX}/\text{\LaTeX}$, if the emitter is extensible, which is guaranteed by the \LaTeX XML binding language: To transform a $\text{\TeX}/\text{\LaTeX}$ document to a given XML format, all \TeX extensions³ must have “ \LaTeX XML bindings” binding, i.e. a directive to the \LaTeX XML emitter that specifies the target representation in XML.

1.3 Generating OMDoc from \gTeX

The \gTeX packages (see Section 2) provide functionalities for marking up the functional structure of mathematical documents, so that the \LaTeX sources contain enough information that can be exported to the OMDoc format (Open Mathematical Documents; see [Kohlhase:OMDoc1.2]). For the actual transformation, we use a \LaTeX XML plugin [LaTeXMLsTeX:github:on] that provides the \LaTeX XML bindings for the \gTeX packages.

1.4 Conclusion

The \gTeX 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

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

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, $\text{\texttt{STeX}}$ can serve as a conceptual interface between the document author and MKM systems: Technically, the semantically preloaded $\text{\texttt{L}^A\text{T}^E\text{T}X}$ documents are transformed into the (usually XML-based) MKM representation formats, but conceptually, the ability to semantically annotate the source document is sufficient.

The $\text{\texttt{STeX}}$ macro packages have been validated together with a case study [Kohlhase04:stex], 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.

1.5 Licensing, Download and Setup

The $\text{\texttt{STeX}}$ packages are licensed under the $\text{\texttt{L}^A\text{T}^E\text{T}X}$ Project Public License [LPPL], 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).

1.5.1 The $\text{\texttt{STeX}}$ Distribution

The $\text{\texttt{STeX}}$ packages and classes are available from the Comprehensive $\text{\texttt{T}^E\text{T}X}$ Archive Network (CTAN [CTAN:on]) and are part of the primary $\text{\texttt{T}^E\text{T}X}$ / $\text{\texttt{L}^A\text{T}^E\text{T}X}$ distributions (e.g. TeXlive [TeXLive:on] and MikTeX [MiKTeX:on]). The development version is on GitHub [sTeX:github:on], 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 $\text{\texttt{T}^E\text{T}X}$ / $\text{\texttt{L}^A\text{T}^E\text{T}X}$ 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 MP
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
```

With this installation using $\text{\texttt{STeX}}$ is as painless as using $\text{\texttt{L}^A\text{T}^E\text{T}X}$, just make sure the $\text{\texttt{STeX}}$ distribution is where `latex` can find it and run `pdflatex` over the main file.

EdN:1

1.5.2 The $\text{\texttt{STeX}}$ Plugin for $\text{\texttt{L}^A\text{\texttt{T}}_E\text{\texttt{X}}ML}$

For the OMDoc transformation of $\text{\texttt{STeX}}$ documents we use a $\text{\texttt{L}^A\text{\texttt{T}}_E\text{\texttt{X}}ML}$ plugin that provides the $\text{\texttt{L}^A\text{\texttt{T}}_E\text{\texttt{X}}ML}$ bindings for the $\text{\texttt{STeX}}$ packages. For installation and setup follow the instructions at [\[LaTeXMLsTeX:github:on\]](#)¹

¹EdNOTE: We are working on a CPAN submission that should make installations painless.

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 T_EX/L^AT_EX workflow. We will group them by the conceptual level they address. Figure 1 gives an overview.

2.1 The sTeX Distribution

The stex package provides stex.sty that just loads all packages below and passes around the package options accordingly and stex-logo.sty that provides the macros \sTeX and \stex that typeset the sTeX logo.

\sTeX
\stex

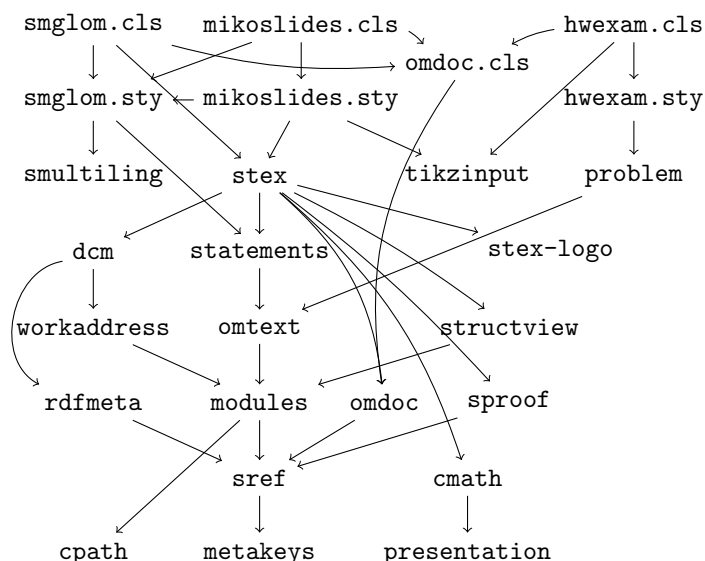


Figure 1: The sTeX packages and their dependencies.

2.2 Content Markup of Mathematical Formulae in T_EX/L^AT_EX

2.2.1 cmath: Building Content Math Representations

The cmath package (see [Kohlhase:cmath:ctan]) supplies an interface for building content math representations. It gives special macros for marking up variables, applications and bindings. It supports the transformation into both OpenMath [BusCapCar:2oms04] and content MathML [CarIon:MathML03].

2.2.2 presentation: Flexible Presentation for Semantic Macros

The `presentation` package (see [Kohlhase:ipsmsl:ctan]) 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.

2.3 Mathematical Statements

2.3.1 statements: Extending Content Macros for Mathematical Notation

The `statements` package (see [Kohlhase:smms:ctan]) 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.

2.3.2 sproof: Extending Content Macros for Mathematical Notation

The `sproof` package (see [Kohlhase:smp:ctan]) 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.

2.3.3 omtext: Mathematical Text

2

2.4 Context Markup for Mathematics

2.4.1 modules: Extending Content Macros for Mathematical Notation

The `modules` package (see [KohAmb:smmssl:ctan]) 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. A side effect of this is that we have an “object-oriented” inheritance mechanism for semantic macros: the semantic macros for the mathematical objects described in a module come with the module itself. As a consequence, the **modules signatures** (only the macro definitions, not the descriptions) need to be loaded before they can be used somewhere else.

²EdNOTE: say something

2.4.2 smultiling: Multilingual Mathematical Modules

In multilingual settings, i.e. where we have multiple \LaTeX documents that are translations of each other, it is better to separate the module signature from the descriptive document.³

2.4.3 structview: Structures and Views

⁴

2.5 Mathematical Document Classes

2.5.1 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`⁵

2.5.2 hwexam: Homeworks and Exams

The `hwexam` package [`Kohlhase:hwexam:ctan`] provides `hwexam.cls` and `hwexam.sty` for marking up homework assignments, and exams. The content markup strategy employed in \LaTeX allows to specify – and profit from – administrative metadata such as time and point counts. This package relies on the `problem` package [`Kohlhase:problem:ctan`] which provides markup for problems, hints, and solutions.

2.5.3 mikoslides: Slides and Course Notes

The `mikoslides` package provides a document class from which we can generate both course slides – via the `beamer` classes – and course notes – via the `omdoc` class – in a transparent way.

2.6 Metadata

2.6.1 rdfmeta: RDFa Metadata for \LaTeX

⁶

2.6.2 dcm: Dublin Core Metadata

⁷

³EdNOTE: continue

⁴EdNOTE: Say something

⁵EdNOTE: continue

⁶EdNOTE: Say something

⁷EdNOTE: Say something

EdN:8

2.6.3 workaddress: Markup for FOAF Metadata

8

2.7 Support for MathHub

The `mathhub` package provides the supplementary packages `mikoslides-mh`, `modules-mh.sty`, `omtext-mh.sty`, `problem-mh.sty`, `smultiling-mh.sty`, `structview-mh.sty`, and `tikzinput-mh.sty` with variants of the user-visible macros that are adapted to the MathHub system – see Section 3.3 for details.

2.8 Auxiliary Packages

2.8.1 metakeys: An extended key/value Interface

EdN:9

9

2.8.2 pathsuris: Managing Relative/Absolute File Paths

EdN:10

10

2.8.3 tikzinput: External TIKZ Pictures as Standalone Images

EdN:11

11

⁸EdNOTE: Say something

⁹EdNOTE: Say something

¹⁰EdNOTE: Say something

¹¹EdNOTE: Say something

3 Workflows and Best Practices

3.1 The “Little Modules” Approach

One of the key advantages of semantic markup with \LaTeX is that the \LaTeX sources are highly reusable by the “object-oriented” inheritance model induced by \LaTeX modules. It turned out to be useful to divide \LaTeX documents into three kinds of files:

1. **module files**: files that essentially contain a collection of \LaTeX modules [KohAmb:smmssl:ctan] – usually a single one whose module name coincides the file name base.
2. **fragment files**: files that contain a group of input references to module- or fragment files – usually one group deep for flexibility, transition text, and additional remarks.
3. **driver files** that set up the document class, contain the preambles, and input reference fragment files.

These correspond to the \LaTeX documents, but can reuse and share \LaTeX fragments and modules. Figure 2 shows a situation, where we have two courses given over multiple years, which results in five course notes documents given by driver files, which share quite a few components. As drivers and fragment files are mostly content-free – they only contribute document structure, this lets all documents contribute from the development of the modules.



Figure 2: Reuse of Fragments and Modules in a Course Notes Setting

The downside of this “object-oriented” inheritance mechanism is that we need to keep the module signatures (see Section 2.4.1) up to date adding to the complexity of document management.

Another advantage of the “little modules approach” is that modules can be developed separately. Indeed all modules of a given subject share common pre- and post-ambles which can be developed in special files – usually named `pre.tex`, `post.tex`, and `preamble.tex` (the latter is included in `pre.tex`). Given such a setup, the call to `pdflatex` can be suitably adapted to handle the pre/postfixes.

3.2 Basic Utilities & Makefiles

The \TeX distribution contains three basic command line utilities to manage \TeX documents in the `bin` directory of the distribution.

`sms` computes the \TeX module signatures for a given \TeX file (see [**KohAmb:smsssl:ctan**] details).

`filedate` and `checksum` that help keeping the metadata of the self-documenting \LaTeX packages in the \TeX distribution up to date.

`installFonts.sh` that installs the fonts necessary for chinese \TeX documents.

These are supplemented by a set of UNIX `Makefiles` in the `lib/make` directory. The way to use them is to include them into a `Makefile` in the directory and then run one of the targets `pdf` and `mpdf` to make the PDF versions of the drivers and modules¹² and `omdoc` and `mods` to generate OMDoc. Note that we need to `make sms` in order to make the respective \TeX module signatures for the modules.

3.3 MathHub: a Portal for Active Mathematical Documents

MathHub (<http://mathhub.info> see [**IanJucKoh:sdm14**]) is a portal for Active Mathematical Documents – documents that are made context-aware and interactive by semantic annotations. \TeX is one of the main input formats for informal active documents. MathHub supports \TeX documents in three ways:

1. MathHub offers free/open hosting in document repositories for (mathematical) \TeX document collections.
2. the backend system supports the large-scale change- and error-management for \TeX documents in the “little modules” paradigm.
3. the front-end displays interactive (HTML5) documents generated from the \TeX sources (via OMDoc).

The MathHub system is probably the best way of developing and hosting larger \TeX document collections. It offers two authoring workflows an online authoring workflow via a direct web interface [**MathHub:oa:on**] or casual users and an offline authoring workflow that we describe next.

3.4 lmh: MathHub’s Build System Locally

As direct web editing workflows are not efficient for larger document collections, the MathHub system offers an offline authoring system. This uses GIT repositories for distribution – the author develops the document collection on a local working copy and then commits for inclusion to MathHub. The MathHub build system can be used locally for efficient development via the `localmh` system [**lmh:github:on**]. In a nutshell – see [**MathHub:law:on**] for details –

¹²EdNOTE: MK: what about the fragments?

1. `localmh` is installed in a docker container that supplies the build system and provides the `lmh` command suite.
2. `lmh pdf` formats \S T E X modules to PDF – building all dependencies, e.g. module signatures, first.
3. `lmh omdoc` generates OMDoc for \S T E X documents – again with dependencies.
4. `lmh xhtml` generates active documents (in XHTML5) from the \S T E X sources or their OMDoc versions.
5. `lmh <gitsc>` distributes the git subcommand `<gitsc>` over multiple repositories.

Various other `lmh` subcommands help with large-scale editing problems like renaming or moving modules, translations in multilingual settings, etc.

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

4.2 The sTeX 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{%
20 \ifundefined{texorpdfstring}%
21 {\let\texorpdfstring\@firstoftwo}%
22 {}%
23 \texorpdfstring{\raisebox{-.5ex}{S\kern-.5ex\TeX}{sTeX}\xspace%
24 }
25 \def\sTeX{\stex}
26 </logo>
```

Index

Numbers written in *italic* refer to the page where the corresponding entry is described; numbers underlined refer to the code line of the definition; numbers in *roman* refer to the code lines where the entry is used.

***, 4, 5 LaTeXML, 5