

System Description

$\mathfrak{S}\text{T}_{\text{E}}\text{X}3$ – A $\text{L}^{\text{A}}\text{T}_{\text{E}}\text{X}$ -based Ecosystem for Semantic/Active Mathematical Documents

Michael Kohlhase, Dennis Müller

Computer Science, FAU Erlangen-Nürnberg

Abstract. We report on $\mathfrak{S}\text{T}_{\text{E}}\text{X}3$ – a complete redesign and reimplementation (using $\text{L}^{\text{A}}\text{T}_{\text{E}}\text{X}3$) from the ground up of the $\mathfrak{S}\text{T}_{\text{E}}\text{X}$ ecosystem for semantic markup of mathematical documents. Specifically, we present: *i)* The $\mathfrak{S}\text{T}_{\text{E}}\text{X}$ package that allows declaring semantic macros and provides a module system for organizing and importing semantic macros using logical identifiers. $\mathfrak{S}\text{T}_{\text{E}}\text{X}3$ is a (now) standard $\text{L}^{\text{A}}\text{T}_{\text{E}}\text{X}$ package with minimal dependencies and is compatible with arbitrary document classes and packages. *ii)* The $\text{R}_{\text{U}}\mathfrak{S}\text{T}_{\text{E}}\text{X}$ system, an implementation of the core $\text{T}_{\text{E}}\text{X}$ -engine in Rust. It allows for converting arbitrary $\text{L}^{\text{A}}\text{T}_{\text{E}}\text{X}$ -documents to XHTML– for $\mathfrak{S}\text{T}_{\text{E}}\text{X}3$ -documents enriched with semantic annotations based on the OMDOC ontology. *iii)* An MMT integration: The $\text{R}_{\text{U}}\mathfrak{S}\text{T}_{\text{E}}\text{X}$ -generated XHTML can be imported and served by the MMT system for further semantic knowledge management services.

This paper uses $\mathfrak{S}\text{T}_{\text{E}}\text{X}3$. The semantically annotated XHTML version of this paper is available at <https://tinyurl.com/cicm22stex>

1 Introduction and History

In the $\mathfrak{S}\text{T}_{\text{E}}\text{X}$ project [sLX], we explore how established communication and publication workflows – this mainly means $\text{L}^{\text{A}}\text{T}_{\text{E}}\text{X}$ in Mathematics and theoretical sciences – can be extended semantically for computer support. The central element of this endeavour is the $\mathfrak{S}\text{T}_{\text{E}}\text{X}$ package [Koh08; sTeX] which allows to *semantically preload* $\text{L}^{\text{A}}\text{T}_{\text{E}}\text{X}$ documents via special (semantic) macros.

$\mathfrak{S}\text{T}_{\text{E}}\text{X}$ documents can be processed by `pdflatex` in the usual way. Additionally, in $\mathfrak{S}\text{T}_{\text{E}}\text{X}1$ they could initially be processed by $\text{L}^{\text{A}}\text{T}_{\text{E}}\text{X}\text{ML}$ [LTX], a $\text{L}^{\text{A}}\text{T}_{\text{E}}\text{X}$ -to-XML transformer, using a dedicated $\mathfrak{S}\text{T}_{\text{E}}\text{X}$ plugin producing OMDOC [Koh06]. Unfortunately, this plugin was elaborate, implementing the OMDOC-specific behaviour via dedicated Perl bindings for the majority of $\mathfrak{S}\text{T}_{\text{E}}\text{X}$ -macros and was correspondingly difficult to maintain. It was also invasive with respect to the $\text{L}^{\text{A}}\text{T}_{\text{E}}\text{X}\text{ML}$ code base and quickly became incompatible with newer versions of $\text{L}^{\text{A}}\text{T}_{\text{E}}\text{X}\text{ML}$. Furthermore, conversion to OMDOC required the usage of dedicated document classes, rendering $\mathfrak{S}\text{T}_{\text{E}}\text{X}$ incompatible with existing and established authoring workflows, and setting up $\mathfrak{S}\text{T}_{\text{E}}\text{X}$ to work in the first place was prohibitively difficult, involving manually changing core parameters of a user’s $\text{T}_{\text{E}}\text{X}$ system.

Nevertheless, the \sTeX package (and associated classes) have been used to produce extensive course materials (3000+ pages of slides and integrated narrative), ca. 2500 exercise/exam problems, and the SMGloM, a multilingual mathematical glossary [SMG], currently containing ≥ 2250 concepts in English (93%), German (71%) and Chinese (11%). This \sTeX -**corpus**, together with the OMDOC format, have informed the development of the \sTeX packages and document model. All \sTeX content is available as *mathematical archives* [Hor+11] on <https://MathHub.info> and can be browsed on <https://mmt.beta.vollki.kwarc.info/:sTeX>.

While the original \sTeX architecture and realization showed that semantic preloading of mathematical documents and the deployment of active documents based on this is possible given enough motivation, the technical/practical problems mentioned above quickly became a showstopper. Not surprisingly, the use of \sTeX never quite gained much traction outside the authors' research group and collaborative projects. Additionally, the continuing development of the MMT system [RK13; MMT] over the last years similarly drove development of its own variant of the OMDOC format and ontology (MMT/OMDOC).

Consequently, we decided to rethink and reimplement \sTeX from the ground up, using \LaTeX3 , with both the problems with \sTeX1 and the developments of the MMT/OMDOC ontology in mind. The result is the \sTeX3 package and system, which we present in this system description (extensive documentation available at [KM]).

Notably, this very document uses \sTeX3 and its module system within the `llncs` document class, and is available (and compiles) on Overleaf at <https://www.overleaf.com/read/tcnwysdzthwx>. We occasionally refer to the source files available there for clarification. In this document, \sTeX was configured such that every semantic macro generates a link to our document server, but it should be noted that this behavior can be fully customized.

2 The \sTeX3 package

The design of the \sTeX3 system was based on the following guiding principles:

1. *Ease of set-up*: The \sTeX3 package should work with a vanilla, unmodified $\text{\TeX}/\text{\LaTeX}$ system – e.g. a sufficiently recent \TeX Live installation – without the need of changing any \TeX -parameters and without any external software.
2. *Universality*: The \sTeX3 package should be compatible with arbitrary \TeX document classes, packages, and authoring workflows. Semantically annotating existing environments (*theorems*, *definitions*, *proofs* etc.) should not impact document layout: their layout should be fully customizable.
3. *No code duplication*: The functionality of \sTeX3 macros and environments should be governed by the \LaTeX -code of the package alone (as opposed to dedicated macro bindings for OMDOC export that implement the same functionality with a different output format) to help maintainability.
4. *MMT-completeness*: \sTeX3 should be a full surface language for MMT/OMDOC.

Let us see how the current system is doing on these accounts.

1. Ease of set-up Indeed, $\S\TeX 3$ now works with any unmodified \TeX system with a $\LaTeX 3$ kernel later than February 2022. For older, but not too outdated $\LaTeX 3$ versions (up to \TeX Live 2018 as running on Overleaf), the missing functionality can be easily added (in this document via the `stex-expl-compat` package).

2. Universality The $\S\TeX 3$ package can be imported in the usual manner (via `\usepackage{stex}`) and only depends on three other packages, namely `ltxcmds`, `standalone` and `xspace`, all of which are ubiquitous, non-invasive and do not take package options that might lead to conflicts.

To allow for collaboration (e.g. via git) and compatibility with submission systems (e.g. arxiv.org) $\S\TeX 3$ can “persist” all semantic macros and other module content into a `.sms`-file during compilation (similar to the `.toc`-file), which can be used in subsequent compilations, obviating the need for the (potentially many) original modules to be physically present. This file can then be put under version control or distributed alongside the document.

To be adaptable to document styles, $\S\TeX 3$ determines the specific highlighting for symbols via four macros, which can be redefined, namely `\compemph{}`, `\symrefemph{}`, `\defemph{}` and `\varemph{}`. For this document, these are defined in `highlights.tex`.

While $\S\TeX 1$ declared its own environments `definition`, `example`, `theorem` etc., doing so necessarily made $\S\TeX 1$ incompatible with document classes that predefine these environments (like `llncs`), or a user’s preferences. However, this was necessary to allow for providing additional semantic information, e.g. as in `\begin{definition}[for=foo,type=inductive]`.

In $\S\TeX 3$, we instead use environments `sdefinition`, `sexample`, `sassertion` etc., that take care of the semantic information provided, but whose typesetting can be customized. For example, by setting `\stexpatchdefinition{\begin{definition}}{\end{definition}}`, every `\begin{sdefinition}[...]` will process the arguments provided, and then delegate to the `definition`-environment for layout and numbering. Analogously, `\stexpatchassertion[theorem]{\begin{theorem}}{\end{theorem}}` will delegate every `\begin{sassertion}[type=theorem,...]` to the `theorem`-environment.

3. No code duplication This principle lead to the following design choice: Rather than converting $\S\TeX$ documents to OMDOC directly, we have the $\S\TeX$ package insert semantic annotations into a non-PDF output format; e.g. XHTML. The package itself determines the full MMT URIs for all symbols, governs the OPENMATH syntax tree and introduces annotations via merely three macros that a “backend” of choice should provide:

- `\stex_annotate:nnn{key}{value}{code}` annotates `code` with `key=value` (e.g. by wrapping `code` in a `...`).
- `\stex_invisible:n{code}` exports `code`, but hides it in the presentation (e.g. by setting `style="display:none"`).

- Lastly, `\begin{stex_annotate_env}{key}{value}` acts like `\stex_annotate:nnn{key}{value}{code}`, but as an environment.

The file `stex-backend-pdflatex.cfg` contains the implementations for these macros for the standard `pdflatex` backend.

4. MMT *completeness* with respect to the \TeX 3 package is a complex issue – at least when trying to avoid code duplication – since the MMT/OMDOC ontology supplies very powerful representational primitives, and will therefore be treated in a regular companion paper submitted to CICM.

3 OMDoc and Mmt

\TeX-XHTML: In \TeX 1, translating \TeX content to OMDOC was achieved directly via \LaTeX XML. In \TeX 3, we instead translate \TeX content to XHTML, augmented with annotations via XML attributes corresponding to the OMDOC ontology. In principle, this workflow allows for a plurality of systems as translators, such as \LaTeX XML or \TeX 4ht. In practice, unfortunately, it has turned out to be difficult to preserve the intended attribute annotations using a current version of \LaTeX XML in math-mode, where they are most important. For now, we therefore implemented our own plain \TeX -interpreter from the ground up using Rust, for converting \LaTeX documents to XHTML. The resulting $\text{Rus}\text{\TeX}$ software¹ uses a user’s local \LaTeX system, keeping the number of required primitives to implement to a reasonable minimum, and can therefore handle (in principle) arbitrary \TeX code to the virtually same degree as the user’s \TeX system (`pdflatex`, to be precise), at the cost of (a priori) no special treatment of higher-level \LaTeX macros (although $\text{Rus}\text{\TeX}$ allows for providing dedicated bindings for these, too). MMT bundles and interfaces with $\text{Rus}\text{\TeX}$ via the JNI to convert to XHTML using MMT’s build system and cache \TeX modules across conversion tasks.

XHTML-MMT/OMDOC: Having obtained semantically annotated XHTML, we have implemented a new XHTML import in MMT to extract the semantic annotations and map them directly to the corresponding MMT/OMDOC concepts. In addition to thus converting \TeX content to OMDOC, the MMT system can host the generated XHTML in a semantically informed manner and offer the full suite of available knowledge management services for \TeX , up to and including type checking and inference.

4 Ongoing and Future Work

Having solved many of the previous problems surrounding \TeX 1 that discouraged users from using \TeX , the most pressing issues now are related to finding,

¹ <https://github.com/slatex/RusTeX>

managing, and reusing *existing* \LaTeX content. We are therefore working on a dedicated IDE for \LaTeX in the form of a Language Server Protocol server and a plugin for VS Code that bundles the MMT system, offers convenient interfaces to interact with it, allows for searching available \LaTeX content (online and locally) and generally helps with semantically annotating documents.²

References

- [Hor+11] Fulya Horozal et al. “Combining Source, Content, Presentation, Narration, and Relational Representation”. In: *Intelligent Computer Mathematics*. Ed. by James Davenport et al. LNAI 6824. Springer Verlag, 2011, pp. 212–227. URL: https://kwarc.info/frabe/Research/HIJKR_dimensions_11.pdf.
- [KM] Michael Kohlhase and Dennis Müller. *The sTeX3 Package Collection*. Tech. rep. URL: <https://github.com/slatex/sTeX/blob/main/doc/stex-doc.pdf> (visited on 04/24/2022).
- [Koh06] Michael Kohlhase. *OMDoc – An open markup format for mathematical documents [Version 1.2]*. LNAI 4180. Springer Verlag, Aug. 2006. URL: <http://omdoc.org/pubs/omdoc1.2.pdf>.
- [Koh08] Michael Kohlhase. “Using \LaTeX as a Semantic Markup Format”. In: *Mathematics in Computer Science 2.2* (2008), pp. 279–304. URL: <https://kwarc.info/kohlhase/papers/mcs08-stex.pdf>.
- [LTX] Bruce Miller. *LaTeXML: A \LaTeX to XML Converter*. URL: <http://dlmf.nist.gov/LaTeXML/> (visited on 03/12/2021).
- [MMT] *MMT – Language and System for the Uniform Representation of Knowledge*. URL: <https://uniformal.github.io/> (visited on 01/15/2019).
- [RK13] Florian Rabe and Michael Kohlhase. “A Scalable Module System”. In: *Information & Computation* 0.230 (2013), pp. 1–54. URL: <https://kwarc.info/frabe/Research/mmt.pdf>.
- [sLX] *sLaTeX: An Ecosystem for Semantically Enhanced \LaTeX* . URL: <https://github.com/sLaTeX> (visited on 03/11/2021).
- [SMG] *SMGloM: A Semantic, Multilingual Terminology for Mathematics*. URL: <http://smglom.mathhub.info> (visited on 04/21/2014).
- [sTeX] *sTeX: A semantic Extension of TeX/LaTeX*. URL: <https://github.com/sLaTeX/sTeX> (visited on 05/11/2020).

² Available at <https://github.com/slatex/sTeX-IDE>