# System Description
# sTeX – A LaTeX-based Ecosystem for Semantic/Active Mathematical Documents

Michael Kohlhase, Dennis Müller, Jan Frederik Schaefer

Computer Science, FAU Erlangen-Nürnberg

## 1  Introduction and History

In the sLaTeX project [sLX], we explore how established communication and publication workflows – this mainly means LaTeX in Mathematics and theoretical sciences – can be extended semantically for computer support. The central element of this endeavour is the sTeX package [Koh08; sTeX] which allows to **semantically preload** LaTeX documents via special (semantic) macros. sTeX documents can be processed by `pdflatex` in the usual way, or via LaTeXML [LTX], a LaTeX-to-XML transformer, which has a sTeX plugin. The semantic annotations are exported into the generated PDF or OMDoc [Koh06] respectively where they can be used for added-value services. The sTeX packages (and 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 [Koh14], currently containing $\geq 2250$ concepts in English (93%), German (71%) and Chinese (11%). This **sTeX corpus** together with the OMDoc/Mmt format have informed the development of the sTeX packages and document model.

While the original sTeX architecture and realization showed that semantic preloading of the mathematical documents and the deployment of active documents based on this is possible given enough motivation, scalability and the management of shared content – one of the potential side-benefits of semantic preloading – quickly became a problem. As a consequence we now host and manage all sTeX content as **mathematical archives** [Hor+11] on `https://MathHub.info` and extended sTeX with special path functionality for cross-references. As a side effect, MathHub can host the interactive HTML generated from the OMDoc in a central location.[1]

In spite of this, the use of sTeX never quite gained much traction outside the authors' research group and collaborative projects. In this system description we detail the effort over the last two years of making sTeX much more usable.

---

[1] e.g. SMGloM under `https://beta.mathhub.info/library/group/c21nbG9t`

## 2 The sTeX EcoSystem

### 2.1 Simplification of sTeX Workflows

Working with sTeX so far required using several external tools and modifying LaTeX parameters, mostly related to sTeX's *module system.*

**Local paths** To find the actual source files containing modularly imported sTeX content, the previous workflow necessitated the creation of a `localpaths.tex` for every top-level `.tex` file, that stores the local file path to sTeX repositories. This workflow has been significantly simplified recently, by replacing the `localpaths.tex` files by a single `MATHHUB` system variable, which points ot the local MathHub clone. LaTeX can now access without needing the `--shellescape` flag.

**SMS mode** sTeX allows the introduction of new semantic macros within `module` environments, using the `\symdef` command. Semantic macros defined in some external module are made locally available with the `\importmodule` command. Since modules can (and in practice often do) contain arbitrary additional content, for `\importmodule` to work, the *semantic* content of a module needs to be extracted from a `module` environment. Previously, this was achieved by a perl script that heuristically parsed a file `foo.tex` for `\symdef` and similar commands, and extracted those into a separate `foo.sms` file. `\importmodule` and related commands then used the `.sms` file to selectively load only the semantic content. This required both an external tool and posed a change management problem - every change to an sTeX module required rerunning the perl script to ensure the `.sms` file is consistent with the source `.tex` file.

The usage of `.sms` files has now been deprecated. Instead, `\importmodule` and related commands enter an "SMS mode" before inputting the required `.tex` file directly, in which everything other than selected sTeX commands (such as `\symdef`) is ignored, obviating the need for external tools or change management considerations. The overhead of multiply reading the narrative content of included files – redundant in SMS mode – turned out to be negligible.

**File stack size** The availability of a module system can quickly lead to deeply nested dependency trees on modules (and hence `.tex`-files), especially when using SMGloM modules. By default, LaTeX has a (relatively low) upper limit for its *file stack*, determining the number of individual files that the TeX engine is allowed to have open at once. Consequently, sTeX users were advised to manually increase the file stack size of their local LaTeX distribution, something most LaTeX users are unfamiliar with or need admin rights.

This has been recently resolved by fully utilizing the difference between *reading* a module (in SMS mode) and *activating* a module. During SMS mode (when the containing `.tex` file is *open*), all semantic macros are merely stored in a separate helper macro. Only after the file has been fully read (and closed by LaTeX), its content is *activated* by executing all semantic macros therein, avoiding recursive `\importmodule`-calls and ensuring that LaTeX's file stack only ever increases

by 1 for semantic imports, obviating the need to manually increasing the file stack size.

**_Standalone SMGloM files_** Previously, sTeX distinguished **documents** (with LaTeX preamble and `document` environment) with **module** (without) and used external build tools to provide modules with preambles on the fly. With the need for external tools – see above – otherwise alleviated, we realized that the LaTeX `standalone` package to make modules independently compilable by `pdflatex`: `standalone.sty` allows for using `\input` on `.tex` files that themselves have a header, without LaTeX throwing an error. This had been possible earlier, but now it is documented best practice.

## 2.2 sTeXLS: An sTeX Language Server and IDE

The highly fragmented structure[2] of sTeX corpora can be a challenge when creating and editing sTeX content. Some of the difficulties can be alleviated with an IDE for sTeX. To avoid being tied to a specific editor, the sTeX IDE is based on **sTeXLS**, a language server that could be used for any editor supporting the language server protocol. sTeXLS has its roots in a bachelor's thesis [Pli18], which explored machine-learning-based approaches to find missing annotations of term references in sTeX documents. The student behind [Pli18] has continued working on sTeXLS, which now supports various features that help authoring sTeX content. Aside from enabling simple interactions like cross-file definition look-up, a key feature of sTeXLS is its ability to point out semantic problems in the source files (**semantic linting**). This ranges from minor issues like redundant imports to actual errors like references to non-existent concepts. sTeXLS addresses such errors by e.g. listing modules from which the concepts could be imported or by suggesting similar sounding concepts in case of a spelling mistake. These features are so useful that sTeXLS is now commonly used for the creation of sTeX content.

## 2.3 Generating Supplemental Material from sTeX Sources

The semantic annotations allow deriving a number of supplemental resources for an sTeX document. Concretely, we have tools to automatically generate dictionaries, glossaries and dependency graphs e.g. to supplement the lecture notes. These tools act directly on the sTeX sources, utilizing the the uniform structure of semantic annotations.

Dictionary generation exploits that definienda in a (natural language) definition are explicitly annotated by a semantic concept. E.g. the synonyms "*linear ordering*" and "*simple ordering*", as well as the German translation, "*lineare Ordnung*", are identified as same concept. An English-to-German dictionary

---

[2] As TeX cannot load document fragments natively, it is natural to prefer very small source files that only contain small semantically self-contained fragments

would then have two entries, one for "*linear ordering*" and one for "*simple ordering*". To generate the dictionary for a lecture, we only include concepts that have been referenced in the lecture notes/slides.

If we map words to their definition rather than their translations, we have a generated glossary. To make it "definitionally closed", we also include entries for all concepts referenced elsewhere in the glossary.

We already routinely generate dictionaries and glossaries for some of our lectures, which was appreciated by the students. We have also experimented with the generation of concept dependency graphs but our visualization efforts had limited success so far due to the sheer size of resulting graphs.

## 2.4  sTEXML2: Partial Preloading and XHTML Harvesting

So far, to obtain formal content from sTEX documents, these documents needed to be converted to OMDOC. To subsequently enable KM services, the resulting OMDOC needs to be additionally converted to the specific OMDOC dialect used by the MMT system by splitting it into (formal) *content* OMDOC and (informal) *narrative* omdoc. This workflow requires:

1. Dedicated sTEX document classes for OMDOC,
2. an sTEX-Plugin for LATEXML that allows generating OMDOC, partially overriding core methods of LATEXML,
3. a suite of LATEXML bindings for most (if not all) sTEX primitive macros, and
4. an MMT component for importing the OMDOC generated by LATEXML.

All of these components needed to be consistently kept in-sync with respect to any updates regarding the sTEX-package, LATEXML, and MMT, and as a result regularly suffered from bitrot and increasingly bloated and difficult to maintain implementations. Additionally, OMDOC generation was incompatible with the document classes commonly used by LATEX users.

As a result, we have deprecated the direct OMDOC generation via LATEXML and the sTEX-Plugin and are re-basing the sTEX packages on a very selective set of semantic primitives. Note that we only need to implement LATEXML bindings for these and can reuse majority of sTEX functionality implemented in TEX – LATEXML covers enough TEX/LATEX primitives by now. This makes it much easier to maintain coherence between the LATEX implementation and the LATEXML bindings. We now use LATEXML to generate XHTML (which LATEXML supports natively) with OMDOC-annotations (provided by the package bindings alone). Crucially, this is compatible with all existing LATEXML bindings, and yields documents that can be immediately inspected with a web browser without loss of document content or significantly impacting layout, maintaining the narrative structure of the original document while introducing *partial* OMDOC information where induced by semantic macros. Afterwards, the MMT system can harvest the generated XHTML to extract the OMDOC fragments relevant for KM services. This simple change of approach realizes an old desideratum in the sLATEX project: flexibly mixing (partial) sTEX functionality into arbitrary LATEX document classes – this is called "sTEX light" in [KKL10].

## 3 Conclusion & Future Work

We have significantly improved the user-friendliness of the sLATEX ecosystem by minimizing the number of *required* external tools and simplifying the general workflow, while supplying additional *optional* tools for added-value services.

As future work, we intend to

1. extend the sTEX language to subsume all primitives of the MMT/OMDOC ontology, allowing sTEX to serve as a full surface language for MMT,
2. allow for semantic markup of arbitrary (in particular informal natural language) document fragments, and
3. improve and extend IDE support, e.g. by providing direct access and search functionality for the SMGloM.

## 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. ISBN: 978-3-642-22672-4. URL: https://kwarc.info/frabe/Research/HIJKR_dimensions_11.pdf.

[KKL10]   Andrea Kohlhase, Michael Kohlhase, and Christoph Lange. "sTeX – A System for Flexible Formalization of Linked Data". In: *Proceedings of the $6^{th}$ International Conference on Semantic Systems (I-Semantics) and the $5^{th}$ International Conference on Pragmatic Web*. Ed. by Adrian Paschke et al. ACM, 2010. ISBN: 978-1-4503-0014-8. DOI: 10.1145/1839707.1839712. arXiv: 1006.4474v1 [cs.SE].

[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.

[Koh14]   Michael Kohlhase. "A Data Model and Encoding for a Semantic, Multilingual Terminology of Mathematics". In: *Intelligent Computer Mathematics 2014*. Ed. by Stephan Watt et al. LNCS 8543. Springer, 2014, pp. 169–183. ISBN: 978-3-319-08433-6. URL: https://kwarc.info/kohlhase/papers/cicm14-smglom.pdf.

[LTX]   Bruce Miller. *LaTeXML: A LATEX to XML Converter*. URL: http://dlmf.nist.gov/LaTeXML/ (visited on 03/12/2021).

[Pli18]   Marian Plivelic. "Using machine learning to support annotating of keywords in mathematical texts". B.Sc. Thesis. FAU Erlangen-Nürnberg, Feb. 2018. URL: https://gl.kwarc.info/supervision/BSc-archive/blob/master/2018/Plivelic_Marian.pdf.

[sLX]   *sLaTeX: An Ecosystem for Semantically Enhanced LATEX*. URL: https://github.com/sLaTeX (visited on 03/11/2021).

[sTeX]   *sTeX: A semantic Extension of TeX/LaTeX*. URL: https://github.com/sLaTeX/sTeX (visited on 05/11/2020).