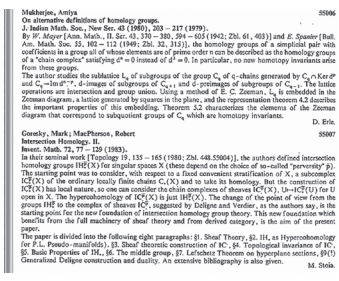
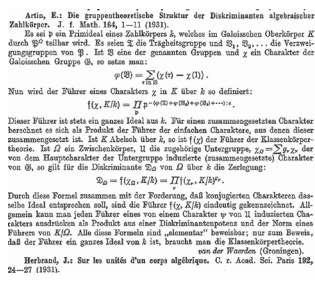


Four Decades of T_EX at zbMATH

Moritz Schubotz (FIZ Karlsruhe, Berlin, Germany) and Olaf Teschke (FIZ Karlsruhe, Berlin, Germany)

In April 2019, zbMATH was completely transformed from T_EX to L^AT_EX sources. On this occasion, we give a brief history of typesetting Zentralblatt volumes, and describe the challenges, methods and benefits of the transition.

A very short overview of typesetting Zentralblatt
T_EX and L^AT_EX have been the standard tools for creating documents for at least two generations of mathematicians. Today it is almost inconceivable that mathematical content was typeset before. Indeed, mathematical typesetting has a much longer history at zbMATH. zbMATH has existed for about 150 years if one includes its printed predecessors, *Jahrbuch über die Fortschritte der Mathematik* and *Zentralblatt für Mathematik*. Although it would be somewhat tempting to retell the history of mathematical typesetting in these periods based on the appearance of the *Jahrbuch* and *Zentralblatt* volumes in a similar manner as, e.g., [5], this would go beyond the scope of the present article. Instead, the figure below provides impressions of volumes based on various lead types, linotypes, or phototypesetting using IBM golfball typewriters.



Formulae from Zentralblatt volumes 1 and 529.

In general, mathematics was always very expensive to typeset.¹ The various developments until the 1970s aimed to make this process more efficient. Famously, the quality of the more cost-efficient phototypesetting technique decreased in comparison to classical lead techniques [5]. Indeed, in the *Zentralblatt* volumes from the years of phototypesetting one comes across several formulae which needed to be added by hand, before one could start to create the phototypesetting master by cutting and glueing.

However, at *Zentralblatt*, the shift to phototypesetting was inevitable due to the growing amount of con-

1 Typically, there were only a few highly specialised typesetters involved. One anecdote that has passed through generations was that one single typesetter was responsible for producing Zentralblatt volumes for many years. Eventually, he was able to spot errors in mathematical formulae without any semantic knowledge.

tent, which resulted in the production of immense register volumes. This created the most obvious and urgent need for further digitisation.

From the beginnings of T_EX to its adaptation

Donald Knuth developed the T_EX system to overcome the limitations of phototypesetting. As a welcome side effect, the fact that T_EX is an open system established the autonomy of mathematical writing. With hindsight, this appears to be an inevitable development, though in the transition period it was not. For instance, a subject of detailed debate was whether it would be more efficient to employ scientists as their typesetters, instead of having specialised staff [9].

A major milestone in this process was the T_EX82 version with both its improvements and stability. The 1982 meeting of the T_EX Users Group at Stanford was the first TUG meeting lasting for about a week. Incidentally, a *Zentralblatt* editor – who was involved at that time in the editors’ exchange program with Mathematical Reviews² – took part, and advocated later in Berlin the adaptation of T_EX for the production of the *Zentralblatt* volumes. At Mathematical Reviews T_EX became fully productive as of 1985, after several years of preparation [7]. In contrast, it took until 1992 before the advent of T_EX at *Zentralblatt*. In the meantime, the phototypesetting technique became no longer sustainable. In 1984, a proprietary, internal Springer system was employed that resembled many of the typesetting functions of T_EX. The commands of the Springer system were designed for specialised technical staff instead of self-use. A key argument for the Springer system in contrast to T_EX was the resulting lower number keystrokes for volume production.

From today’s viewpoint, the greatest advantage of moving away from phototypesetting was that at least the texts and formulae were available in a digital form for the volumes 531–734. However, the disadvantage was that after the Springer system became outdated in the early 1990s, the introduction of T_EX required conversion: simultaneously to the Springer-based production plan. Since the transition to T_EX also included a migration from Springer servers to FIZ Karlsruhe, the expected delay occurred.³ Fortunately, the T_EX expertise acquired in the meantime allowed for a successful transition. It turned out that even most formulae could be translated automatically, though some constructions

² That was also the time when a merger of both services was actively pursued.
³ The 3-month hiatus between Vol. 734 and 735 is by far the largest post-war gap.

had inherent problems; e.g., matrices needed conversion from column-based encoding to row-based encoding.⁴

From T_EX to L^AT_EX

For some years after the switch to T_EX, the production of printed volumes had been the main objective. The main objective was the appropriate presentation. A standardised encoding, which is desirable from an information retrieval viewpoint, was less relevant. This changed in the second half of the 1990s after the online database became the primary objective. In 2004, adapting to these needs, the database production switched to a PostgreSQL system. This stored all available information in ASCII-coded T_EX texts. This framework hadn't changed significantly until recently.

Of course, despite the impressive robustness and stability of T_EX, technical development didn't stop in the 80s. Probably already in the mid-90s L^AT_EX was the preferred dialect for many users. Today L^AT_EX accounts for more than 90% of zbMATH review submissions. Since the database routines previously interacted with T_EX, reviews in L^AT_EX required a re-standardisation. Furthermore, an increasing amount of data is provided in the UTF-8 format. Its conversion to T_EX encoding for internal storage and later reconversion for online presentation may cause errors and a loss of information. This pertains especially to bibliographic reference data which may include the need to encode native script such as Arabic, Chinese, Farsi, Hebrew, Japanese, Korean, or Russian. This makes it desirable to have the X_YL^AT_EX expansion available.

How does one convert ~20 million formulae?

Therefore, a switch of the production system underlying zbMATH to X_YL^AT_EX/L^AT_EX was in preparation for several years. The initial work-intensive step was made when MathML was introduced in 2010 [1]. Standard tools for MathML conversion employed by zbMATH, such as Tralics [8], require L^AT_EX source. Thus, it was necessary to convert different T_EX commands to L^AT_EX – at least, those commands which could be processed by MathML converters. This part of the conversion could be amply addressed by regular expressions. However, it turned out to be necessary to build the full expression tree for T_EX formulae. This step, which was finally taken in April of this year, was preceded by an upgrade of PostgreSQL, which allows for UTF-8 handling and a 14-hour routine that converted approximately 18 million standard inline mathematical expressions and 600,000 displaystyle formulae. These included more sophisticated environments, like `\alignat` or `\gathered`. Fortunately, only a few environments, e.g., commutative diagrams, require additional manual transition. Overall, the introduction of L^AT_EX resulted in a pause of zbMATH updates of about a month. In contrast to the introduction of T_EX, which took three months, this is a significant improvement. The most vis-

ible difference for zbMATH users looking at the review sources is the replacement of `$....$` by `\(...\)`. Standardisation will allow for much easier integrity checks, and for a much more seamless integration of the submitted reviews. Additionally, considerably better presentation is possible by new functions available via L^AT_EX packages. MathJax (needed for maths presentation in browsers not capable of MathML) works better on widely used L^AT_EX commands instead of less frequent T_EX commands. Another advantage pertains to the MathML generation: while Tralics works well for presentation purposes, it hasn't been developed further for some years, so it is reasonable to look for alternative options like L^AT_EXML [6]. The availability of L^AT_EX code makes such alternative implementations much more feasible.

Paving ground for future developments

Even more importantly, further developments in mathematical information retrieval and processing will most likely be based on L^AT_EX. L^AT_EX became the de facto standard in mathematics, and most working mathematicians use L^AT_EX in their publication workflows. Moreover, most websites use L^AT_EX as an input language for mathematical formulae. To make mathematical content better discoverable, multiple approaches exist for enhancing semantics in mathematical formulae. For example, the NIST Digital Library of Mathematical Functions developed a set of semantic LaTeX macros for mathematics. These macros are easy to use for mathematicians fluent with L^AT_EX. By using these macros, with minimal overhead, information retrieval systems would be able to better disambiguate the syntax and semantics for mathematical expressions. Eventually, this will provide better search and recommendation functionality for users of mathematical digital libraries [2]. To some extent, such approaches have already been applied to realise the zbMATH formula search [4, 3], but having standardised L^AT_EX sources at hand will certainly make further developments in this direction much more feasible.

References

- [1] P. Baier, O. Teschke, Zentralblatt MATHMLized, *Eur. Math. Soc. Newsl.* 76, 55–57 (2010; Zbl 1278.68088)
- [2] H.S. Cohl et al., Growing the digital repository of mathematical formulae with generic LaTeX sources, *Lect. Notes Comput. Sci.* 9150, 280–287 (2015; Zbl 06512424)
- [3] F. Müller, O. Teschke, Full text formula search in zbMATH, *Eur. Math. Soc. Newsl.* 102, 51 (2016; Zbl 1366.68355)
- [4] M. Kohlhase, H. Mihaljević-Brandt, W. Sperber, O. Teschke, Mathematical formula search, *Eur. Math. Soc. Newsl.* 89, 56–58 (2013; Zbl 1310.68217)
- [5] D.E. Knuth, Mathematical typography, *Bull. Am. Math. Soc., New Ser.* 1, 337–372 (1979; Zbl 0404.92025)
- [6] L^AT_EXML, a L^AT_EX to XML/HTML/MathML converter, <https://dlmf.nist.gov/LaTeXML/>
- [7] D.C. Latterner, W.B. Wolf, T_EX at *Mathematical Reviews, Tugboat* 10, No. 4, 639–654 (1989)
- [8] Tralics: a L^AT_EX to XML translator, <https://www.sop.inria.fr/marelle/tralics/>
- [9] Letters et alia, *Tugboat* 4, No. 2, 81–102 (1983), <https://www.tug.org/TUGboat/tb04-2/tb08letters.pdf>

⁴ By looking at the sources of these, a reader can easily indicate that they were not genuine T_EX-coded; e.g., single variables were not set in formula italics in the old system, so they lack conversion until now.



Moritz Schubotz [moritz.schubotz@fiz-karlsruhe.de] is a senior researcher for mathematical information retrieval and open science. He maintains the support for mathematical formulæ in Wikipedia and is off-site collaborator at NIST.

Olaf Teschke's photo and CV of the author can be found in previous Newsletter issues.