

# BIB<sub>TE</sub>X<sub>ML</sub>: An XML Representation of BIB<sub>TE</sub>X

Luca Previtali, Brenno Lurati, Erik Wilde

ETH Zürich (Swiss Federal Institute of Technology), Switzerland  
luca@bibtexml.org, brenno@bibtexml.org, dret@bibtexml.org

## Abstract

BIB<sub>TE</sub>X<sub>ML</sub> is an XML representation of BIB<sub>TE</sub>X data. It can be used to represent bibliographic data in XML. The advantage of BIB<sub>TE</sub>X<sub>ML</sub> over BIB<sub>TE</sub>X's native syntax is that it can be easily managed using standard XML tools (in particular, XSLT style sheets), while native BIB<sub>TE</sub>X data can only be manipulated using specialized tools.

**Keywords:** <sub>TE</sub>X, <sub>La</sub><sub>TE</sub>X, BIB<sub>TE</sub>X, XML, bibliography

## 1 Introduction

<sub>La</sub><sub>TE</sub>X [3] is an advanced text formatter which is widely used for typesetting in mathematics, physics and engineering. BIB<sub>TE</sub>X [4] is a program and file format designed by PATASHNIK and LAMPORT in 1985 which provides an easy way to manage bibliographic references for <sub>La</sub><sub>TE</sub>X. The BIB<sub>TE</sub>X format is character and field (tag) based.

Many BIB<sub>TE</sub>X bibliographies are available online. These collections contain huge amounts of references (for an example see the size of <http://liinwww.ira.uka.de/bibliography> shown in Table 1). The main problem of these collections is that BIB<sub>TE</sub>X is a simple format, which does not allow complex queries and data manipulation.

References	> 1.1 million
Amount of data	660 MBytes
Cross references	16'000
Number of URLs	100'000

Table 1: Example of bibliographic database size

The goal of BIB<sub>TE</sub>X<sub>ML</sub> is to develop an XML environment for representing and structuring BIB<sub>TE</sub>X bibliographies, which makes the management of bibliographic data easier, and to build an online database which allows upload and download of bibliographic entries.

## 2 Proposed Solution

XML provides an efficient way to structure data using an easy to edit and readable format. With XSLT it is possible to easily convert XML data into various formats.

The basic idea is that the whole bibliography should be managed using XML. `.bib` files are no longer edited: they are generated only for BIB<sub>TE</sub>X use. To enable a reuse of all existing BIB<sub>TE</sub>X collections, we need a conversion tool for transforming BIB<sub>TE</sub>X bibliography into BIB<sub>TE</sub>X<sub>ML</sub> format (as described in Section 3.2). Once the BIB<sub>TE</sub>X<sub>ML</sub> representation is obtained, it is possible to store the BIB<sub>TE</sub>X<sub>ML</sub> entries in an online database (see Section 4). This database provides complex queries and data navigation which help the user to fetch the required references (see Figure 1).

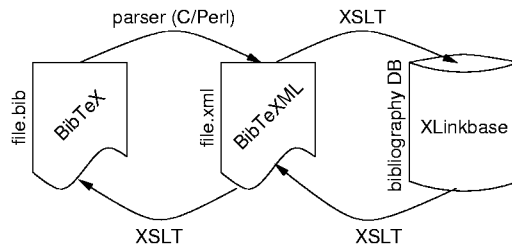


Figure 1: BIB<sub>TE</sub>X<sub>ML</sub> architecture

## 3 Implementation

### 3.1 BIB<sub>TE</sub>X to BIB<sub>TE</sub>X<sub>ML</sub> Parser

BIB<sub>TE</sub>X does not have a formally defined grammar, which makes the implementation of a parser difficult. BEEBE [1], who works on BIB<sub>TE</sub>X tools since 1990, created a prototype BIB<sub>TE</sub>X grammar based on the available documentation and various experimental tests. Our parser is built on top of the *bibparse* tool (a *lex/yacc* generated lexical analyzer) developed by BEEBE.

BIB<sub>TE</sub>X entries may contain <sub>La</sub><sub>TE</sub>X-coded special characters or commands (eg, `\'e`, `\LaTeX{}`) which are (in this <sub>La</sub><sub>TE</sub>X-representation) not meaningful in XML. To solve this problem, we implemented a *Perl* converter which translates Unicode compliant characters into XML entities (eg, `\'e`  $\rightarrow$  `&#x00E8;`) and inserts a special XML element for not Unicode-coded strings (eg, `\LaTeX{}`  $\rightarrow$  `<tex code="\LaTeX{ }">LaTeX</tex>`). This solution provides full <sub>La</sub><sub>TE</sub>X backwards compatibility and the possibility to have an XML-compliant representation.

### 3.2 BIB<sub>TE</sub>X<sub>ML</sub> Format

BIB<sub>TE</sub>X ignores unknown fields, which makes it possible to extend the set of usable tags. BIB<sub>TE</sub>X<sub>ML</sub> provides two different *XML schemas*. The first one defines all standard BIB<sub>TE</sub>X fields (eg, *author*, *title*, *editor*, ...), while the second defines non-standard extensions (eg, *ISBN*, *URL*, *abstract*, *language*, ...). This allows us to expand the BIB<sub>TE</sub>X<sub>ML</sub> capabilities maintaining full compatibility with the original format.

A typical BIB<sub>TE</sub>X entry looks like this:

```
@Book{lampport:86,  
  author = "Leslie Lamport",  
  title = "{\LaTeX}: A Document  
  Preparation System",  
  publisher = "Addison-Wesley",  
  year = "1986"  
}
```

After the conversion, the BIB<sub>T</sub>E<sub>X</sub>ML entry has the following format:

```
<bibliography>
  <bibitem type="book">
    <label>lamport:86</label>
    <author>
      <firstname>Leslie</firstname>
      <lastname>Lampport</lastname>
    </author>
    <title><tex code="\LaTeX{">LaTeX</tex>:
      A Document Preparation System</title>
    <publisher>Addison-Wesley</publisher>
    <year>1986</year>
  </bibitem>
</bibliography>
```

BIB<sub>T</sub>E<sub>X</sub>ML provides a macro environment helping to reduce the amount of stored data and make it more manageable. For example, each author may be stored only once, and individual author entries may reference it, enabling a simple implementation of complex queries and easy navigation within the database. Furthermore, every feature of BIB<sub>T</sub>E<sub>X</sub>'s macro mechanism is kept in the BIB<sub>T</sub>E<sub>X</sub>ML macro format. Additionally, it is possible to define more structured macros as shown in the following example:

```
<macro id="450">
  <firstname>Leslie</firstname>
  <lastname>Lampport</lastname>
  <email>lamport@pa.dec.com</email>
</macro>
...
<author id="450">
...

```

Macros can be referenced with a simple attribute inside the element.

### 3.3 BIB<sub>T</sub>E<sub>X</sub>ML to BIB<sub>T</sub>E<sub>X</sub> Conversion

BIB<sub>T</sub>E<sub>X</sub> processes BIB<sub>T</sub>E<sub>X</sub>-formatted entries only, and it is therefore necessary to convert BIB<sub>T</sub>E<sub>X</sub>ML data into the BIB<sub>T</sub>E<sub>X</sub> format. The XML structured bibliography is easily transformed into a `.bib` file using XSLT (see Figures 1 and 2). This can be locally done using the BIB<sub>T</sub>E<sub>X</sub>ML2BIB<sub>T</sub>E<sub>X</sub> XSLT style sheet. Thanks to the XML environment, it is also possible to easily translate BIB<sub>T</sub>E<sub>X</sub>ML entries into various formats (eg, HTML or plain text) using easily adaptable XSLT style sheets.

### 3.4 XLinkbase

While BIB<sub>T</sub>E<sub>X</sub>ML is intended to be used as an exchange format and a way to represent BIB<sub>T</sub>E<sub>X</sub> in XML, we see it as an intermediary format for our *XLinkbase* system, which is designed for managing large amounts of highly interlinked information, using a data model similar to *Topic Maps* [2]. XLinkbase makes it possible to easily browse and manipulate BIB<sub>T</sub>E<sub>X</sub> entries, while BIB<sub>T</sub>E<sub>X</sub>ML is used as import and export format for the system.

Figure 2 shows the overall model of interaction, where users are working on the XLinkbase data, and only export BIB<sub>T</sub>E<sub>X</sub>ML if it is required for processing with the BIB<sub>T</sub>E<sub>X</sub> program, or for exchange or other purposes (such as generating an HTML list of a number of bibliographic entries).

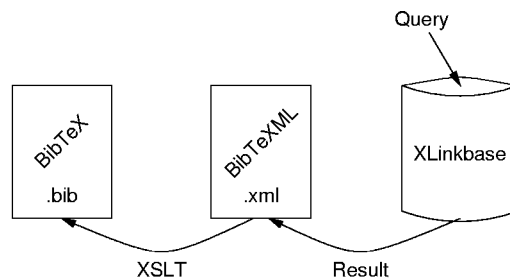


Figure 2: XLinkbase and the BIB<sub>T</sub>E<sub>X</sub>ML format

## 4 bibtexml.org

On <http://bibtexml.org/> it is possible to use and test some of the presented utilities and find more detailed information regarding the structure and the implementation of the BIB<sub>T</sub>E<sub>X</sub>ML project. In particular, it offers the following services:

- Check the correctness of a BIB<sub>T</sub>E<sub>X</sub> file.
- Convert BIB<sub>T</sub>E<sub>X</sub> to BIB<sub>T</sub>E<sub>X</sub>ML and vice versa.
- Download the conversion utilities.
- Get the BIB<sub>T</sub>E<sub>X</sub>ML XML Schema definitions and XML Namespace definitions.
- Browse through a sample BIB<sub>T</sub>E<sub>X</sub>ML XLinkbase.

We are also working on setting up a repository of BIB<sub>T</sub>E<sub>X</sub>ML bibliographic entries in the near future, which can then be searched using a subset of XPath expressions.

## 5 Conclusions

In this poster we describe BIB<sub>T</sub>E<sub>X</sub>ML, an XML syntax for BIB<sub>T</sub>E<sub>X</sub>ML data. The not formally defined BIB<sub>T</sub>E<sub>X</sub> grammar does not allow an efficient data manipulation. Formatting bibliographies entries with XML enable us to obtain a clean and effective structure, thus making possible the creation of more powerful bibliographic databases and their manipulation using general-purpose XML-tools. The conversion tools that we have developed allow the translation of available BIB<sub>T</sub>E<sub>X</sub> collections into the new format, and the BIB<sub>T</sub>E<sub>X</sub>ML database can also be translated to the original BIB<sub>T</sub>E<sub>X</sub>-compatible format.

## 6 Acknowledgments

We would especially like to thank Nelson Beebe for the valuable feedback.

## References

- [1] NELSON F. H. BEEBE. Bibliography Prettyprinting and Syntax Checking. *TUGboat*, 14(4):395–419, December 1993.
- [2] INTERNATIONAL ORGANIZATION FOR STANDARDIZATION. Information technology – SGML Applications – Topic Maps. ISO/IEC 13250, 2000.
- [3] LESLIE LAMPOR. *L<sup>A</sup>T<sub>E</sub>X: A Document Preparation System*. Addison-Wesley, Reading, Massachusetts, 1985.
- [4] OREN PATASHNIK. Bib<sub>T</sub>E<sub>X</sub>ing. Technical report, February 1988.