

SGML and T_EX in Scientific Publishing

N.A.F.M. Poppelier
Elsevier Science Publishers
Physical Sciences and Engineering Division
R&D Department
email: n.poppelier@elsevier.nl

Abstract

Elsevier Science Publishers has for a few years investigated the possibility of accepting compuscripts, a manuscript in electronic form, created with T_EX, L^AT_EX and a few other text processing systems, and converting these to SGML form. This paper will discuss the current status of these activities, the reasons for converting compuscripts to SGML form, and the various ways in which T_EX is used.

Introduction

Until a few years ago the results of scientific research were always published in the following way: the author writes an article using whatever method he prefers. Some authors use T_EX, some authors a word processor, others prefer a typewriter, and there are still authors who consider the pen to be the best desktop publishing tool that is available.

When the author is satisfied with his work, he or she submits the manuscript, say, to the editor of a journal, usually in duplicate or triplicate. The editor enters the manuscripts into his administration, and sends a copy to one or more "referees," who review and judge the manuscript. This is one of the places on the long road from manuscript to published paper where the publisher adds something to the paper. The system of peer reviewing, which provides a kind of scientific certification, is one example of "added value." The quality of the printed product and the world-wide distribution of an article in a journal can also be labelled "added value."

The referee(s) can either recommend acceptance of the manuscript for publication, possibly with suggestions for alterations, or rejection of the manuscript. If it is accepted, and after the author has put it into final form, the manuscript is sent to the publisher, where it is adapted to the style of the journal by a technical editor. Then the manuscript is re-typed (typeset) and, in most cases, stored in electronic form. Because of the fact that the material has been re-typed, it may contain errors. These errors are removed by proof-reading, which is a task that can be performed, for example, by the technical editor or the author.

Finally, the paper is published, together with other papers, in the journal and after a while the author sees his article in print.

The time that passes between submission of the manuscript and publication of the article in printed form can vary between a few weeks and a year, or even longer. In some fields of science this period of time is too long and authors sometimes ask publishers why this time cannot be reduced.

Modern Times

The conventional method of publishing, which I have sketched above, can probably not be made any more efficient by using conventional means. However, many authors use their own text processing systems to produce their papers and reports. This text can be transmitted in electronic form. In principle, this makes it possible to

- publish the paper within a shorter length of time;
- extract bibliographic information and abstracts, and store this information in a database.

In the first case we see that, by making use of manuscripts in electronic form, which I will call *compuscripts* in the rest of this paper, that is by making use of modern electronic means, the conventional method of publishing scientific journals can be made more efficient. Furthermore, the production of secondary publications, such as abstract journals, can be made more efficient since the abstracts do not have to be re-typed.

An important consequence is that the publisher can now add extra value to the journal or to his entire range of scientific publications, by making the

material available in various electronic forms, such as hypertext books, databases on CD-ROM, etc.

However, a publisher can only accept a compuscript if the instructions of the text processing system of the author can, in one way or another, be converted to that of the publisher, which in most cases is a computer-driven phototypesetter. Since a publisher wants to work as efficiently as possible, this conversion must be automatic or nearly automatic. This means that the compuscript must be prepared by the author according to a set of rules. In most cases, of course, a technical editor still has to read a draft printed on paper to mark appropriate changes to the text, the typography or the layout. Especially if the paper contains tables or mathematical formulae, this is an intricate task.

Research at ESP

Elsevier Science Publishers (ESP) has for a few years been investigating the possibility of accepting compuscripts created with modern text processing systems. The aim of our present research is to develop a method for manipulating the compuscripts in such a way that the material can be published more efficiently than before, and is also stored for re-use at a later stage.

We take the view that the transition from manuscript-based to compuscript-based publishing must at least preserve the flexibility of the current submission procedures and, where possible, improve upon it. This means, for instance, that the author should not be bothered with compuscript preparation instructions that are significantly different between publishers, or even between the various journals of one publisher. It also means that publishers should be able to accept compuscripts produced by a range of text processing systems, and submitted either by electronic network or on a diskette, via ordinary mail.

SGML

As Figure 1 shows, the Standard Generalized Markup Language (SGML) [1, 2, 3] plays a central role in our approach. SGML promises to become a very important tool, since it allows us to (i) separate the logical structure of the document from its visual structure, i.e. its appearance; (ii) treat all accepted compuscripts, once they have been converted to SGML form, in a uniform working environment; (iii) handle the compuscripts independently of output medium and output format; (iv) store the text, wholly or partially, in a database and use the database contents again at a later stage for various

other forms of publications; (v) standardize compuscript input formats between the various scientific publishers, thus helping to preserve the flexibility of submission mentioned above.

It is not our intention to ask authors to submit SGML-coded compuscripts yet. However, we foresee that, as a consequence of the fact that the departments of the US Government as well as major industries have adopted SGML as a standard for document interchange, developers of text processing software will provide their customers with facilities for conversion from their particular text format to SGML. Authors will then be able to produce their compuscripts with the text processing system of their choice and to provide the publisher with an SGML version without additional effort.

Until then, in order to have all material submitted in the form of compuscripts available in SGML-coded form, conversion facilities must be developed. The scheme we would like to realize is indicated in Figure 1.

On the one hand, Elsevier Science Publishers wants to allow authors to use their favorite text processing system. On the other hand we know that there are dozens of text processing systems that are being used by authors in various disciplines.

This means that we have to make a choice. Two factors influence this choice: (i) the text processing systems chosen must correspond to a fair proportion of the total number of authors, and (ii) a compuscript prepared with a certain text processing system must, ideally, be convertible to SGML form. The first point is something that should be investigated separately for every discipline. The second point implies that we must choose text processing systems that, whenever possible, produce texts where the logical structure is, at least partially, indicated. In other words: a Word compuscript based on a style sheet is to be preferred over a plain Word file, and \LaTeX is preferred over plain \TeX .

\TeX Compuscripts

\TeX will be one of the text processing systems included in our electronic publishing scheme. Later this year, as a first step towards having a \TeX compuscript scheme for all appropriate journals, *ESP* will start accepting \LaTeX -coded compuscripts for a few of their physics journals. For various reasons we have chosen \LaTeX as the most convenient variety of \TeX .

- Since according to Lamport [4], “the primary function of almost all the \LaTeX commands ... [is] to describe the logical structure of [a] docu-

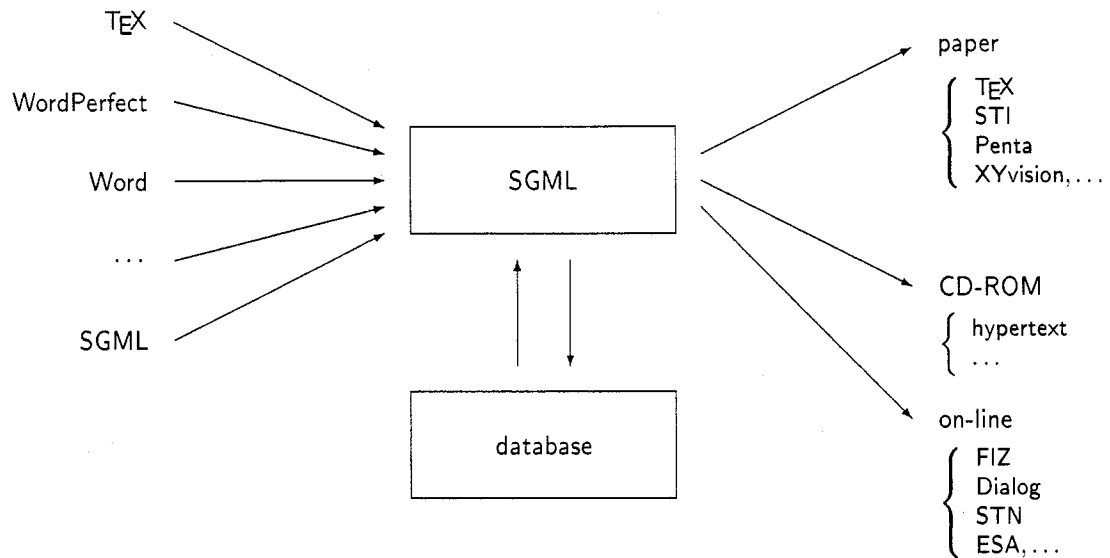


Figure 1: Conversion to and from SGML.

ment” conversion of a L^AT_EX-coded compuscript to SGML form is possible.

- The idea of a document style, which is fundamental to L^AT_EX, enables us to produce camera ready copy for a particular journal by changing only the `\documentstyle` command in the submitted compuscript.
- L^AT_EX is relatively easy to learn.
- L^AT_EX is described in a book, so that it can be considered as a *de facto* standard.
- Especially for large review articles and for books, the following L^AT_EX tools come in very handy
 - cross-referencing with symbolic keys,
 - automatically generated table of contents,
 - index and bibliography tools.

Status Quo

So far we have achieved several things. First of all, we have created L^AT_EX document styles for four of ESP’s scientific journals. These new document styles exactly reproduce the layout of the corresponding journals, which so far are still produced in the conventional way.

We have also been able to develop a set of programs for complete automatic conversion of L^AT_EX documents with mathematical formulae to SGML form. To be more precise: we convert L^AT_EX documents with mathematical formulae to a *document type definition* [1, 2] for scientific papers that is based on the one developed by the *Association of American Publishers (AAP)*.

A document type definition is a description of the logical structure of a certain class of documents, using a method that resembles the specification of the syntax of a programming language. As an example, the document type definition of a novel is represented in a diagrammatic manner in Figure 2.

Conversion of L^AT_EX-coded tabular material to SGML form, i.e, to the AAP document type definition or possibly a different document definition, is currently under investigation. Another automatic conversion that is currently under investigation is the conversion from Word, with an appropriate style sheet, to SGML.

Book Projects with TeX

Recently the Physical Sciences and Engineering Division of ESP has published several books that have been written in plain TeX or in L^AT_EX, and the number of TeX-coded and L^AT_EX-coded books that we will publish in the near future will probably increase. Our experience is that L^AT_EX is an excellent tool for these larger projects due to the separation of logical structure and visual structure, the ability to do cross-referencing with symbolic keys and the presence of tools for index and bibliography; these advantages were already mentioned under the section *TeX Compuscripts*.

However, we have also noticed that making a TeX-coded or L^AT_EX-coded book ready for publication requires more time when the authors do not have to follow a set of instructions, but are free to type things the way they think is best.

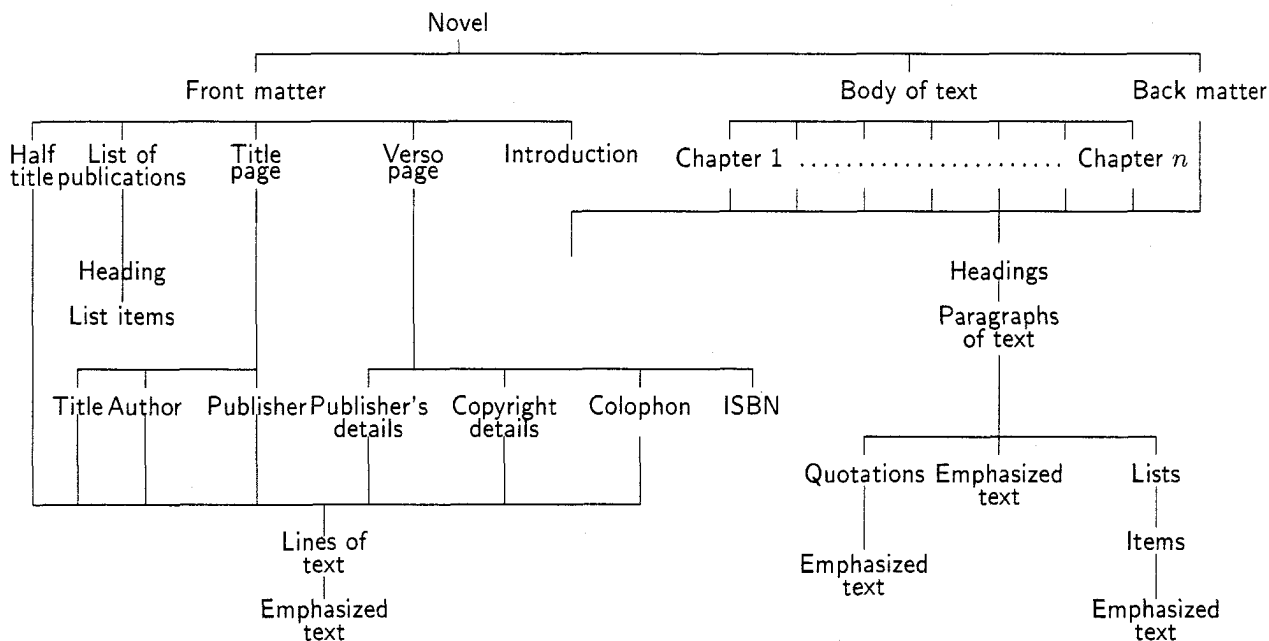


Figure 2: Document type definition of a novel.

Other Uses of T_EX

Lots of programs, such as database programs, have output capabilities that are somewhat meager. After all, a database program is meant to store data, not to print it in a sophisticated manner. The fact that T_EX is a programmable text formatting system makes it an excellent “print engine” for such programs. As long as the database program is capable of inserting some fairly simple texts in its output – most database programs have report generators that are able to accomplish this – T_EX’s capabilities as a text formatter can be used to print the output in any desired way.

An example of this approach is symbolic mathematics packages that can compute complicated formulae and generate the T_EX codes for formatting them. Another example is provided by the database-publishing activities of the *Excerpta Medica Publishing Group (EMPG)*, which is part of the Biomedical Division of *ESP*.

The *EMPG* employs a database of article openings, abstracts and citation lists that have been derived from articles in a selected set of biomedical journals. The information in this database is stored in a structured form. With this information the *EMPG* produces abstract journals that are called the *Core Journals in <...>*, where <...> is, for instance, *Ophthalmology*, *Cardiology* or *Neurol-*

ogy. The *EMPG* has been experimenting with using L^AT_EX for the production of *Core Journals*: the relevant portions of the database are extracted and converted to L^AT_EX. The L^AT_EX-coded form of this information can then, in principle, be used to produce the camera ready copy for the journal, using a document style that was developed in-house.

In the near future, starting this year, a similar system will be implemented for the article openings of all journals that are published by *ESP*. The purpose of this project, which is called *CAPCAS*, is to have the article openings, i.e. title, author(s), abstract, keywords and publication history, available in SGML-coded form and to store this information in a database. From this database we can then create secondary publications of various kinds, such as volume indexes, author indexes and abstract journals. Also for this type of database-publishing, T_EX is used to format the structured information on paper.

Conclusion

In this paper I have sketched the ideas that *Elsevier Science Publishers* have developed concerning the handling of author prepared manuscripts in electronic form. In this scheme, SGML will play a central role, and T_EX, in one or more of its varieties, will also play a part.

Furthermore, I have given a few examples of other ways in which T_EX's text formatting capabilities can be put to good use, namely in the fields of book production and database publishing.

Summarizing, I think it's no exaggeration to say that, because of its programmability, T_EX is eminently suitable for large-scale text production, both directly as a text processing system, and indirectly as the output component of a database-publishing system. Considering that

- the quality of material typeset in T_EX is considered satisfactory even by very demanding users,
- T_EX runs on an impressive number of different types of computers, and
- T_EX output can be printed on a wide range of output devices,

the number of uses of T_EX will undoubtedly increase in the years to come.

Bibliography

- [1] International Standard ISO 8879: *Standard Generalized Markup Language (SGML)*. ISO (1986).
- [2] Martin Bryan: *SGML, an author's guide to the Standard Generalized Markup Language*. Addison Wesley, (Workingham, 1988)
- [3] Eric van Herwijnen: *Practical SGML*. Kluwer Academic, (Dordrecht, 1990).
- [4] Leslie Lamport: *L^AT_EX, a document preparation system*. Addison Wesley, (Reading MA, 1986).