

BIJLAGE EE**SGML (, T_EX and . . .)****C.G. van der Laan**

September 1990

Abstract

What SGML (and T_EX) is all about is given in a nutshell. Markup of example document elements, by SGML and L^AT_EX, are provided. Coupling SGML to T_EX is considered by direct translation and by the intermediate procedural markup phase. Interfacing SGML to (La)T_EX is also addressed. Some guidelines are provided in order to decide when SGML, or T_EX (alone, both, or neither) might be beneficial. It is a 3-in-1 paper: what is SGML and T_EX all about, examples of marked up copy in SGML and (La)T_EX and the coupling issues, finished up with a literature compilation.

What is a Document?

The Association of American Publishers (AAP) *Reference Manual on Electronic Manuscript Preparation and Markup* defines a document as:

A document is an organized collection of smaller pieces of text (such as chapters) and images (such as figures) that are called elements. The elements in a document have a relationship to each other which gives the document a definite organization called document structure.

Lifecycle-phases of documents

Manuscript preparation requires that additional information be interspersed within the text to aid any subsequent processing. That information, called markup, is usually specific to a particular publisher, system or printing device. A universal standard method of marking up electronic manuscripts, however, offers many advantages.

The *complete Lifecycle* of a document can be thought of as:

- preparation,
- distribution,
- reading,
- storing (Paper? Electronically? Optically?),
- other usage, reuse?¹

Standard Generalized Markup Language (SGML) supports the complete Lifecycle, where *future* usage of the document is not necessarily restricted to printing. SGML must be complemented, however, by generally accepted Document Type Definitions (DTDs). The Association of American Publishers [AAP,1987 and 1989] and the

British National Bibliography [Smith, 1987] have provided some DTDs. In order to serve the primary aim of publishing coupling to formatters must be supported too.

T_EX supports formatting and electronic document exchange.

What is SGML?

SGML provides us a language to describe documents. SGML has it made possible to achieve two goals:

1. establish a standard means of identifying and tagging parts of an electronic manuscript so that computers can differentiate between these parts; and
2. provide some logical ways of representing special characters, symbols and tabular material, using only the ASCII character set usually found on standard keyboards.

SGML is defined in ISO8879 [1986]. An introduction for SGML is given by Barron [1989], a gentle introduction is included in Sperberg-McQueen & Burnard [1990]. Introductory books are Bryan [1988] and van Herwijnen [1990].

Purpose

The purpose of SGML is to facilitate *information exchange*, and reusability (in other contexts, even yet unknown contexts),

—*Then and There*—

via a description *language*, where information is packed in documents, containing, text, graphics, formulas, tables, etc.

¹We can talk about reuse and rework. From the author's point of view we are dealing with rework. Publishers like to reuse copy.

Meta Language

SGML is a *meta* language, which can be used to define an arbitrary number of markup languages in a standardized way. This means, for any class of documents, markup rules can be prescribed by SGML, yielding a *language*, the Document Type Definition, for that class. The parser checks compliance of the marked up copy to the DTD.

Standard

Formerly: no consensus on markup “codes”
(WordPerfect, Wordstar, MacWrite, ...; Scribe, troff, T_EX, L^AT_EX, ...)

Presently: SGML ISO standard

Standard ^{def} It can be used to define an arbitrary number of markup languages in a *standardized* way.

Entails:

general applicability,
longer longevity,
improved reusability,
enhanced exchange possibilities.

Generalized

Formerly: (typeset) *marks* for *specific* “here and now” printers, via *direct* markup.

Presently: Marks are *generic*.² This is done with procedural markup. Macro calls are inserted as markup tags, where the implementation of the macros (the format or style file) represents the style, accounts for the fonts, etc. The printer hardware is shielded by intermediate languages. Intermediate language copy is printed via drivers. Change of style needs another style file, no modification of tagged copy is necessary. Change of printer hardware needs another driver, no modification of tagged copy nor modification of format file!

Generalized ^{def} Abstraction from the specific (printing) to the general (other usage), by emphasizing the *structure* of a document and to specify intent without regard for appearance.

Markup

Formerly: (typeset) *marks* in the margin (Marks are bound to a version; no “data-integrity”)

Presently: Marks are integrated along with copy (Data-integrity of markup code is preserved.)

Markup ^{def} Term used to describe codes added to the electronically prepared document.

Author’s point of view

Authors have to markup their copy with

1. awareness of the DTD which applies to the document; either the DTD must be understood or templates must be available;
2. knowledge of which (begin) tags to use where and how;
3. knowledge of ranking, attribute use, tag minimization;
4. knowledge about how to create entity references.

These aspects are treated in author’s guidelines. The above can be alleviated by providing an SGML *environment*, or better, a document preparation environment, supported by menus and templates with prompts. I agree with the general expectation that authors can concentrate on structure and content by using a standard (generic) markup language, or a sufficiently advanced document workbench, leaving formatting issues to publishers, or software vendors. Because of this “separation of concerns” the author’s task is simplified.

Publisher’s point of view

Publishers make use of sufficiently accepted DTDs. They provide authors with guidelines and proof tools. DTD writing requires knowledge of the various types of markup, such as presentational, direct, procedural and descriptive markup.

Example markups

No markup

In order to remind the unpleasant look of documents with just words, we start with a no markup example

```
TeX A system for formatting text
TeX and the accompanying macro
package LaTeX provide powerful means ...
```

Presentational markup

Documents with tabs, indentations, and in general positional control make use of what is called presentational mark up, in order to convey the meaning

```
TeX:
A system for formatting text.
```

```
TeX and its accompanying macro
package LaTeX provide
powerful means of formatting
text to be output on either
- a simple matrix printer,
- a laser printer or
- a photo typesetter.
```

Presentational markup is functional with poetry, such as Alice’s mousetail as mentioned by Malcolm Clark [1989] or D_EK’s favourite poem of Piet Hein [*The Errors of T_EX*, 1989], where the words are arranged along an ellipse.

²Not specific to print/plot/phototypesetter hardware.

Direct markup

When specific print instructions are included, we get direct markup:

```
@T:                []
A system for formatting text.[]
                        [I]
@T and its accompanying macro
package @LT provide
powerful means of formatting
text to be output
on either                [I]
- a simple matrix printer, [I]
- a laser printer or      [I]
- a photo typesetter.
```

[I] is a print instruction indicating to go to the next line and *indent*; @<name> stands for a process with special format effect.

Procedural (\LaTeX) markup

A markup command, where the implementation of the command contains print instructions, is considered a procedural markup command; when the printer is changed the implementation has to be changed too, not the marked up copy. \LaTeX markup of the example reads

```
\subsection*{\TeX}
A system for formatting text.
\par
  \TeX\ and its accompanying macro
  package \LaTeX\ provide
powerful means of formatting text
to be output on either
\begin{itemize}
\item simple matrix printer,
\item a laser printer or
\item a photo typesetter.
\end{itemize}
```

Descriptive (SGML) markup

Descriptive markup goes even further and uses markup which describes the structure and intent of the various parts of the document:

```
<h>&TeX;
<p>A system for formatting text.
<p>&TeX; and its accompanying macro
  package &LaTeX; provide
  powerful means of formatting
  text to be output on either
<li>
<it>simple matrix printer,
<it>a laser printer or
<it>a phototypesetter.
</li>
```

Formatting information with SGML?

It is possible to convey formatting information via SGML. This is done with element attributes or with Processing Instructions. Other symbols than those in the ASCII character set are often denoted by an entity reference to the font containing those symbols, with appropriate loading of the font elsewhere. With respect to attributes, one can think of specifying open space in order to include illustrations from other (electronic) sources. Also, indication of a representation choice is possible if properly accounted for in the DTD. Consider for example the representation of labels of list items: alphabetical or roman/arabic numeral.

```
<li number=alpha>
<it>a simple printer,
<it>a laser printer or
<it>a photo typesetter.
</li>
```

It is possible in SGML to include document parts which already contain format information. The parser must be told to lay back. For a notation to be allowed it must be declared via e.g.

```
<!NOTATION TeX SYSTEM>
```

for \TeX formatted copy. Appropriate entity and attribute specifications are also needed in the DTD. For authors the equation formatting with the type attribute (value= \TeX) has to be supplied as:

```
<eqn type=TeX>
  $$X\cap(A\cup B)=(X\cup A)\cap(X\cup B)$$
</eqn>
```

Processing Instructions (PIs) can be used to tell the local system how it should process data contained within a document. For example, SETM typography markup instructions:

```
<p><?[s24][sec][rm]>T<?[pri][rm]his>
  paragraph ...
```

In this case the SETM instructions are in brackets, preceded by the PI open delimiter <?. The meaning of this instruction is to treat “T” as “24pt” initial letter to be set using the **roman** version of the face currently defined in the **secondary** type family.

Availability

As mentioned by Herwijnen [1990], Sobemap and The Publisher are some SGML systems that are already available.

Support

Support for SGML is done by the companies, as part of automation projects. There also exists a Dutch chapter of the SGML Users Group.³

³SGML-Holland secretary: D. van Wijnen, Wolters Kluwer. P.O. Box 989, 3300AZ Dordrecht. 078-334933; e-mail: surf003@kub.nl.

SGML User's Group secretary: S. G. Downie, Softquad Inc, 720 Spadina Avenue, Toronto, Ontario M5S 2T9, Canada.

Courses

Courses are also provided by private companies, and the National Normalization Institutes.

What is SGML not?

SGML is not

- WYSIWYG (pronounced wīšewīg, and stands for what you see is what you get),
- a formatter, certainly not a standard formatter.

What is T_EX?

T_EX stands for $\tau\epsilon\chi$, the first three letters of the Greek word for *technique*, which also means art. T_EX is a *machine independent* formatting language designed by Don Knuth, [*The T_EXbook*, 1990 (Version 3.0)]. Michael Doob gives an easy start to T_EX in *A Gentle Introduction to T_EX* [1989]. There also is an introduction in French by Seroul [1989] and various German introductions by Apelt [1988], Schwarz [1989]. Von Bechtolsheim [1990] is impressive. L^AT_EX, by Leslie Lamport [1985], is a macro collection for simplified use of T_EX. L^AT_EX uses *procedural* markup. Buerger [1990] gives a L^AT_EX introduction. Bruin [1989] gives a Dutch introduction to L^AT_EX.

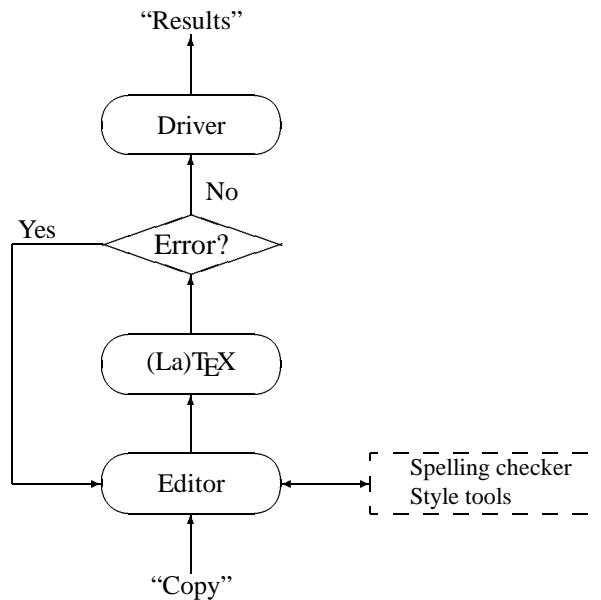


Figure 1: Correction cyclus

Purpose

The purpose of T_EX is “making beautiful books.”

Processing (L^A)T_EX

L^AT_EX is processed in three steps: edit the copy, format the copy to create a dvi file, and print the resultant dvi file. A diagram of this looks like:

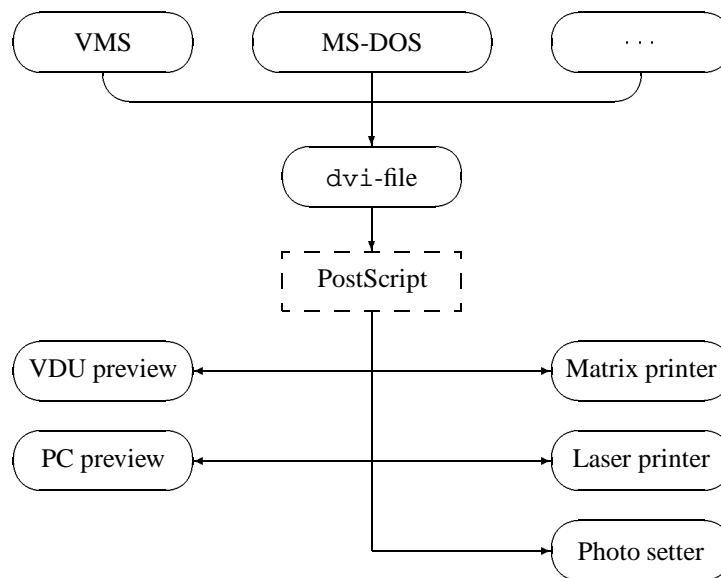
copy $\xrightarrow{\text{editor}}$ ASCII $\xrightarrow{\text{(La)T}_E\text{X}}$ dvi-file $\xrightarrow{\text{driver}}$ results

The more steps to the process, the more cumbersome is correction handling because of larger “loops.” This is the case when the use of T_EX is combined with SGML markup. The SGML parsing and linking extends the loop.

Availability

T_EX is available on many computers under various operating systems with a variety of drivers for previewing (such as VDU), printing, and phototypesetting. Documents written in T_EX or L^AT_EX can be ported easily. Exchanging documents via e-mail is also generally possible except for the incorporated graphics. When graphics are part of the document, T_EX can be combined with Postscript, which is used within the T_EX community. T_EX is in the *public domain*. Drivers are not. They do, however, generally have added value from the companies you buy the driver from. See the ads in any *TUGboat*.

Moreover, T_EX systems can make use of fonts from various sources, such as Adobe’s PostScript fonts and, of course, METAFONT.

Figure 2: (La)T_EX's use

Support

Support is organized by the various users groups. Software, style files, macros etc. are distributed (via e-mail, ftp or floppy disks) by the T_EX Users Group (TUG)⁴; in the Netherlands it is distributed by NTG⁵; in France by GUTenberg; in Germany by DANTE; in the United Kingdom by ukT_EXug; in the Nordic countries by “Nordic TuG”. There is also the new TUGlib service in Utah.

Courses

Courses are organized by TUG and other T_EX user groups, especially in conjunction with their main meetings.

Relationship: DTDs, SGML, T_EX, formats and . . .

The relationship of T_EX, SGML and other applications is illustrated in figure 3. An integrated⁶ implementation is Arbortext's “The Publisher,” which has AAP's DTD s built-in, and requires SUN hardware.

Note that the two backarrows denote some of the work in progress by Elsevier Science Publisher, Bleeker [1989], Poppelier [1990].

⁴Editorial and TUG address: T_EX Users Group, P.O. Box 9506, Providence RI 02940, USA. email: TUGboat@Math.AMS.com.

⁵NTG: Nederlandstalige T_EX Gebruikersgroep. Secretary: G.J.H. van Nes, Postbus 394, 1740AJ, Schagen, 02246-4185; e-mail: vannes@ecn.nl, ntg@hearn.

⁶Ikons user interface, SGML layer, T_EX layer, Postscript handling (optionally); with SGML, T_EX and dvi files as intermediate results

⁷The meta-ness is a strong point —flexibility, adaptability and openness— but, surprisingly at the same time it acts as a weak point —everybody writes or modifies DTD, with discrepancy as result.

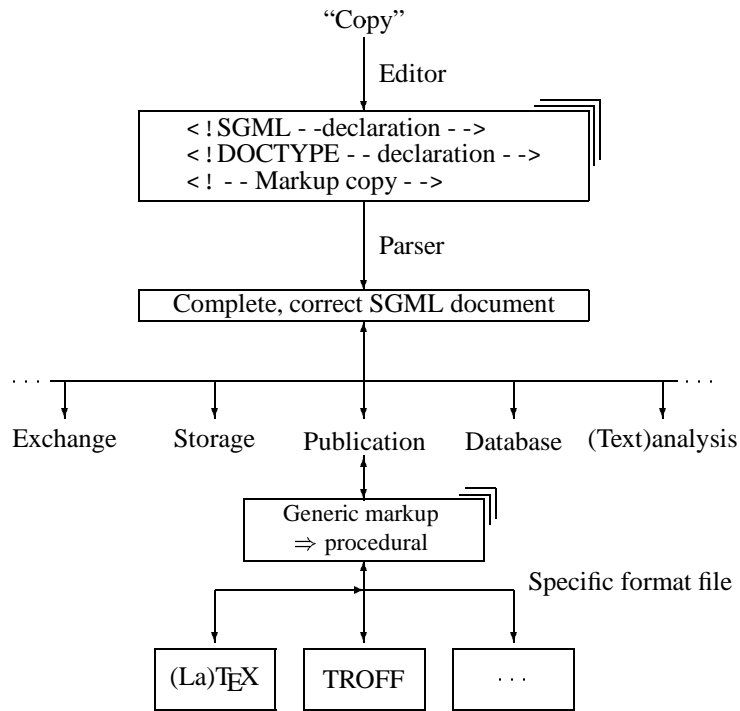
SGML ór T_EX sufficient?

NO, needed are format files and DTDs as well! If a manuscript is printed with T_EX for personal use only, there is nothing to worry about. When no reuse of a document is in sight, but remote publishing and electronic exchange are the case, it pays to use standard format/style files—which reflect the lay-out of the document type—along with T_EX. When reuse or abstract structuring are being used, standard DTDs will be crucial.⁷

Interfacing vs. transformation

Interfacing copy marked up by any formatter to SGML can be prescribed in SGML via the NOTATION mechanism. Of course it has to be incorporated into the DTD and appropriately implemented: the parses should lay back and leave the formatting to the T_EX marked up copy to T_EX.

Transformation SGML into T_EX is different. Using T_EX as a back-end formatter to SGML can be done. It is supported by the link mechanism of SGML. Needed is at least a table of corresponding notations in order to substitute the markup tags from SGML into T_EX. The other way round has to be done by a separate program. In the sequel we will study example document elements marked up by SGML and (La)T_EX; transformation issues will be addressed as well.

Figure 3: Relation SGML and (La)T_EX

Examples

Simple text

As an example let us take the simple text given earlier. The (basic) SGML markup and L^AT_EX markup have been given in previous sections. Coupling comes down to a change of representation, except for the omitted endtags. This direct approach needs the substitutions:

SGML	⇒ T _E X
<h>	\section{
<p>(first)	}
<p>	\par or blank line
	\begin{itemize}
	\end{itemize}
<it>	\item
&TeX;	\TeX
&LaTeX;	\LaTeX

This suggests systematic coupling of all entities and tags to equivalents in L^AT_EX. The explicit endtags are more natural to handle than the omitted ones. The handling of the first and latter occurrences of <p> have to account also for respectively finishing the heading </h> and ending a paragraph </p>, which have been omitted.

Letter

A typical letter has the structure:

- Background
 - Heading (Logo, address, phone, ...)
 - Footer (numbering, ...)

- Context (running heads next pages, ...)
- Reference
- Your reference
- Date
- Addressee (name, company, address, zip code)
- Beginning (Dear...)
- Contents
- End matter (Salutation, name, position)
- Additions (PS, enclosure, cc)

SGML markup

The SGML markup for a typical letter might look something like:

```

<!DOCTYPE letter PUBLIC
-- DTD to be used --
"-//NTG//DTD Letter//EN">
<letter -- start-tag -->
<ref> CGL/Ba/B89-007
<yourref> MC/L1/L89-001
<date> 4 august 1989
<address> Malcolm Clark, ICRF
<email>
    malcolm@icrf.ac.uk
</email>
<dear>Malcolm
<p> Thank you very much ...
...
<p> Some details about the course ...
...
<signed name=CGL>
</letter -- end-tag -->
  
```

L^AT_EX specification

The same letter using L^AT_EX markup might look like:

```
\documentstyle[12pt]{letter}
\address{% return address
  C. G. van der Laan  \\
  \ldots}
\signature{Kees}
\begin{document}
\begin{letter}{% address
  Malcolm Clark  \\
  \dots}
% no ref or your ref
% date is handled automatically
\opening{Dear Malcolm}
\par
Thank you very much \ldots
\begin{quote}
$\vdots$
\end{quote}
Some details about the course \ldots
\begin{quote}
$\vdots$
\end{quote}
\closing{Best regards}% Handles signature
%ps, cc, enclosure all possible
\end{letter}
\end{document}
```

Letter result

Because a sample L^AT_EX letter could not be processed simultaneously in this paper, the printed letter has been omitted.

Transformation

What comes to mind when looking at both representations of marked up copy is the difference in the sequence order of tagged items in SGML and L^AT_EX. With *complete* marked up SGML copy to be transformed into procedural T_EX, there is no problem: in T_EX macros strings can be stored for later usage. This will be shown along with the tabular example in the section transformation revisited.

Bridge card deal

In *TUGboat* 11#2 I have described typesetting bridge using T_EX.

SGML markup

The SGML markup for Figure 4 might look like:

```
<deal><vuln>N/None
  <comm>Deal: demo
<hand n><spades>J74
  <hearts>AJ
  <diams> QJT2
  <clubs> Q874
<hand e><spades>K86
  <hearts>T9542
  <diams> 874
```

```
<clubs> T3
<hand s><spades>QT952
  <hearts>Q83
  <diams> AK5
  <clubs> A6
<hand w><spades>A3
  <hearts>K76
  <diams> 963
  <clubs> KJ952
</deal>
```

(L^A)T_EX specification

The procedural T_EX markup for Figure 4 might look like:

```
\crdima{N/None}{\vtop{\hbox{Deal:}
  \hbox{demo}}}%
{\hand{J74}{AJ}{QJT2}{Q874}}%N
{\hand{K86}{T9542}{874}{T3}}%E
{\hand{QT952}{Q83}{AK5}{A6}}%S
{\hand{A3}{K76}{963}{KJ952}}%W
```

Transformation

The transformation comes down to

SGML	⇒ T _E X
<deal>	\crdima
<vuln>N/None	{N/None}
<comm>DEAL: demo	a suitable \vtop
<hand x>	{\hand

and all the cards per colour surrounded by curly braces, with an extra “}” after the clubs. Although once again a simple example, the translation table is not natural.

T_EX macros

Figure 4 is created by using the T_EX macros:

```
\def\hand#1#2#3#4{%
%Example: \hand{AKJ765}{AK9}{--}{T983}
\vtop{\hbox{\strut\s\enspace#1}
\hbox{\strut\h\enspace#2}
\hbox{\strut\d\enspace#3}
\hbox{\strut\c\enspace#4}}%end \vtop
}%end \hand
%
\def\crdima#1#2#3#4#5#6{%
%purpose: layout bridge hand
%#1 left upper text
%#2 right upper text
%#3, #4, #5, #6: N, E, S, W hands
\vbox{\halign{
&##\quad\cr
#1& #3& #2\cr
$\vcenter{#6}$&$\vcenter{\copy\NESW}$&
&$\vcenter{#4}$\cr
& #5& \cr
}%end \halign
}%end \vbox
}%end \crdima
%
\def\NESWfig{%
\vbox{\font\small=cmr9
\def\str{\vrule height2.2ex%
```

N/None ♠ A3 ♥ K76 ♦ 963 ♣ KJ952	♠ J74 ♥ AJ ♦ QJT2 ♣ Q874 <div style="border: 1px solid black; padding: 5px; width: fit-content; margin: 0 auto;"> N W E S </div> ♠ QT952 ♥ Q83 ♦ AK5 ♣ A6	Deal: demo ♠ K86 ♥ T9542 ♦ 874 ♣ T3
---	--	--

Figure 4: Bridge deal

```

depth.75ex width 0pt}
\offinterlineskip\tabskip0pt\hrule
\halign{\vrule\hskip2pt\relax
##\hfil\tabskip3pt& \str\hfil##\hfil&
##\hskip2pt\relax\hfil\vrule
& \hbox to 0pt{\hss\N\hss}& \cr
\W& \phantom{N}& \E\cr
& \str\hbox to 0pt{\hss\S\hss}& \cr
} %end \halign
\hrule} %end \vbox
} % end \NESWfig
\setbox\NESW\hbox{\NESWfig}

```

SGML requirements

The following DTD is needed:

```

<!ENTITY % ISOpub PUBLIC
"ISO 8879-1986//ENTITIES Publishing//EN">
<!ELEMENT deal - - (vuln, comm?, hand+)>
<!ELEMENT (vuln|comm) - O CDATA>
<!ELEMENT hand - O (spades,
hearts, diams, clubs)>
<!ATTLIST hand nesw (n|e|s|w) #REQUIRED>
<!ELEMENT (spades|hearts|diams|clubs)
- O CDATA>

```

Note. In the DTD we could have imposed sequence ordering by changing `hand+` into `(handn, hande, hands, handw)`. But this requires that all the hands are needed and that is too restrictive.

Some math

The following examples of mathematical formulas are borrowed from the "Mathematical Formulas" report [van der Laan, Coleman, Luyten, 1989]. In this report, SGML and \LaTeX markup are supplied for formulas from various fields: elementary mathematics, set theory, geometry, functional analysis, calculus (differential equations, special functions, continued fraction), statistics, algebra (tensor calculus), homology (diagrams) and quantum mechanics. A few simple ones are selected here.

\LaTeX results

The following was formatted with \LaTeX markup:

$$X \cap (A \cup B) = (X \cup A) \cap (X \cup B)$$

$$x \notin A \not\subset B$$

$$\|a(x+y)\| \leq |a| \cdot (\|x\| + \|y\|)$$

$$\int \frac{1}{\sqrt{1+x^2}} dx = \log(1 + \sqrt{1+x^2})$$

(Basic) SGML markup

To accomplish this with SGML markup, you might enter:

```

<fd>X&cap;(A&cup;B) =
(X&cup;A)&cap;(X&cup;B)</fd>
<fd>x&notin;A&notsubset;B</fd>
<fd><fen d>a(x+y)<rp d</fen>&le;
|a|.(<fen d>x<rp d</fen>
+<fen d>y<rp d</fen>)
</fd>
<fd><in><opd><fr>1</><rad>1+
x<sup>2</rad></fr>dx</in>=
<rf>/log/(1+<rad>1+x<sup>2</rad>)
</fd>

```

The DTD used is an adapted version of AAP's DTD by D.C. Coleman.

(Direct) $T_{E}X$ markup

\LaTeX and $T_{E}X$ markup are very similar for these examples:

$$X \setminus \text{cap} (A \setminus \text{cup} B) = (X \setminus \text{cup} A) \setminus \text{cap} (X \setminus \text{cup} B)$$

$$x \setminus \text{notin} A \setminus \text{not} \setminus \text{subset} B$$

$$\|a(x+y)\| \setminus \leq |a| \cdot (\|x\| + \|y\|)$$

$$\int \text{bfr} 1 \setminus \sqrt{1+x^2} \setminus \text{efr} dx = \setminus \log(1 + \setminus \sqrt{1+x^2})$$

Transformation

To accomplish the SGML to T_EX transformation, some general substitutions are needed:

```
SGML    ⇒ TEX
<fd>    \[ or $$
</fd>   \] or $$
<sup/2/ ~2
etc.
```

For the first set theory example the following substitutions are additionally needed

```
SGML ⇒ TEX
&cap; \cap
&cup; \cup
```

Functional analysis required moreover

```
SGML    ⇒ TEX
<fen d> \ |
<rp d>  \ |
&le;    \leq
```

For the integral the following definition (format) is needed for the fraction, where use is made of </> as parameter separator:

```
\def\bfr#1</>#2\efr{{#1\over#2}}
```

Also needed are the substitutions

```
SGML    ⇒ TEX
<in>    \int
<opd>   \relax
<fr>    \bfr
</fr>   \efr
<rad>   \sqrt{
</rad>  }
<rf/log/ \log
```

Note. A translation table is once again not straightforward; unnatural are the SGML difference in norm open and closing, and the fancy use of </>, i.e. null endtag for numerator and omitted opening tag for denominator. The short reference for the modulus sign is neat.

(Complete) SGML markup

The sobemap parser yielded the following (visually edited) result for the set theory example:

```
<FD DCN="GEO.FORM">
<FL>
X&cap;<FEN STYLE="S" LP="PAR">
  A&cup;B
  <RP STYLE="S" POST="PAR"></FEN>
=
<FEN STYLE="S" LP="PAR">
  X&cup;A
<RP STYLE="S" POST="PAR"></FEN>
  &cap;
<FEN STYLE="S" LP="PAR">
  X&cup;B
<RP STYLE="S" POST="PAR"></FEN>
</FL>
</FD>
```

This shows that complete SGML is verbose. For example, consider the complete tags for parenthesis “(” and “)”. Thanks to the short reference mechanism the input can look natural. Coupling of the above to T_EX can be done. How to handle automatically attributes in the best way is not yet clear to me. It is not efficient for parenthesis, “(” and “)”, to be expanded first by the parser into a fence tag with the appropriate attribute value, followed by substitution at the T_EX level into “(” and “)” again.

Because matrices are treated similarly to tabular material, we have omitted a matrix example, and refer to the next sections, where a table is studied.

AAP's simple table

A simple table is characterized by: simple table entries, one header row, no header subrows, no footer, no intra referencing, and no caption. From the SGML technical point of view no attributes are used. AAP's example simple table [Markup of Tabular Material, 1989], even more simplified, is reproduced in figure 5.

(Basic) SGML markup

The SGML markup for Figure 5, with a minor structural adaptation and some layout modifications, reads:

```
<tbl>
<no>Table AAP
<tt>Job Changes: 1973&ndash;1980
<th>
<th>Gain/Loss of Hospitals since 1973
<th>Total No. of CEO Job Changes
  1973&ndash;80
<th>Survival Rate of CEO's
<bdy>
@Texas|20||22%
@Maryland|5|42|24%
<ft>
<au>David Kinzer
<atl>Turnover Of Hospital Chief
  Executive Officers: A Hospital
  Association Perspective
<nme>Hospital and health Services
  Administration
<dt>May&ndash;June 1982
</tbl>
```

A DTD for simple tables is not separately provided by AAP; it is incorporated as part of the complex table DTD. The “simple table”-example obeys the following SGML structure description

```
<!ENTITY row STARTTAG "row" >
<!ENTITY column STARTTAG "c" >
<!ELEMENT tbl - - (hd, bdy, ft) >
<!ELEMENT hd - O (no?, tt?, th+) >
```

Table AAP Job Changes: 1973–1980			
	Gain/Loss of Hospitals since 1973	Total No. of CEO Job Changes 1973–80	Survival Rate of CEO's
Texas	+20	—	22%
Maryland	+ 5	42	24%

Source: David Kinzer, "Turnover Of Hospital Chief Executive Officers: A Hospital Association Perspective," *Hospital and health Services Administration* May–June 1982.

Figure 5: AAP's simplified table

```

<!ELEMENT bdy      - O row+      >      \hbox to .5\entrywidth{\hss#}\hfil\cr
<!ELEMENT row      - O c+        >      %preamble line
<!ELEMENT ft       - O (au|src|atl|nme|dt) >      \tablerule\noalign{\vskip1ex}
<!ELEMENT          - O (th, c, au, src, atl, nme, dt)-- >      \omit{\bf Table AAP} Job Changes:
          - O (%t.cs;) -- Character string-->      1973--1980
<!SHORTREF tablemap "@ " row      >      \hidewidth\cr
          " | " column      >      \tablerule\noalign{\vskip1ex}
Note. Some tags are presumed part of the general DTD, >      \omit &
e.g. no, tt. >      \omit\vtop{\noindent\hsize=\entrywidth
>      Gain/Loss\nl
>      of Hospitals \nl
>      since 1973}&
>      \omit\vtop{\noindent\hsize=\entrywidth
>      Total No. \nl
>      of CEO \nl
>      Job Changes \nl
>      1973--80} &
>      \omit\vtop{\noindent\hsize=\entrywidth
>      Survival \nl
>      Rate of \nl
>      CEO's} \cr%end header row
>      \noalign{\vskip.5ex\hrule\vskip.5ex}
>      %head-body separation
>      Texas & $+ $20& --- & 22%\cr
>      Maryland& $+ $5&42 & 24%\cr
>      \noalign{\vskip1ex}%body-foot separation
>      \noalign{Source: David Kinzer,
>      ``Turnover Of Hospital Chief Executive
>      Officers:
>      A Hospital Association Perspective,``
>      {\it Hospital and health Services
>      Administration\ /} May--June 1982.
>      }%end \noalign
>      }%end \halign
>      }%end \vbox

\newdimen\entrywidth%\entrywidth=<default>
\newdimen\tablewidth
\tablewidth=.5\hsize\default
\def\nl{\par\noindent}
\def\ndash{--}
\def\ablerule{\noalign{\hrule}}
\newdimen\digitwidth\setbox0=\hbox{\rm0}
\digitwidth=wd0
%?-command for nonsignificant leading
% zeroes, Knuth p241
\catcode'\?=active
\def?{\kern\digitwidth}
%\btbl %AAP's simple example with direct
% TeX markup
\entrywidth=2cm
\tablewidth=4\entrywidth
\vbox{\hsize=\tablewidth
\halign{\hbox to\entrywidth{\# \hss}\hfil&

```

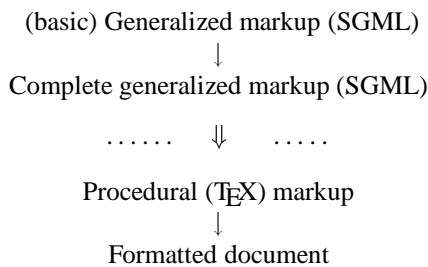
Transformation revisited

Complete SGML markup with procedural T_EX markup representation can be achieved via the SGML link mechanism, or by any automatic substitution process (programmable editor, preprocessor). Translation into T_EX can be done "direct" or via procedural markup. For the

mathematical examples given in van der Laan, Coleman [1989], this has been done by Grootenhuis [1990]. He has used the SGML link mechanism for “substitution” and did direct coupling to L^AT_EX. The *direct* coupling approach, without procedural (L^AT_EX) markup, has the disadvantage that change of formatter requires (T_EX) source adaptation. In order to abstract from any particular format, we have considered procedural T_EX markup as an intermediate phase, in van der Laan, Coleman [1990].

SGML ⇒ T_EX using procedural markup

The “generalized markup ⇒ format” process can be characterized by the following four levels,



with the input phase and output (print) phase before and after. The interfacing with procedural markup is illustrated below, with AAP’s simplified table as example. This four level process resembles the coupling of higher level programming languages such as PASCAL and ADA to FORTRAN (numerical libraries). For more on the latter, I refer to Einarsson [1988] and other early work of mine [1984].

Note. Of course, one can also work the other way round: start from procedural marked up copy and couple that to SGML.

(Completely tagged) SGML markup, AAP’s table

The complete SGML markup version—complete means expansion of short references into tags and addition of omitted (end) tags—of AAP’s simple table reads:

```

<tbl>
<no>Table AAP</no>
<tt>Job Changes: 1973&ndash;1980</tt>
<hd>
  <th></th>
  <th>Gain/Loss of Hospitals
    since 1973</th>
  <th>Total No. of CEO Job Changes
    1973&ndash;80</th>
  <th>Survival Rate of CEO’s</th>
</hd>
<bdy>
<row><tsb> Texas</tsb>
  <c>20</c><c></c><c>22%</c>
</row>
<row><tsb>Maryland</tsb>
  <c>5</c><c>42</c><c>24%</c>
</row>
    
```

```

</bdy>
</tbl>
<au>David Kinzer</au>
<atl>Turnover Of Hospital Chief
  Executive Officers: A Hospital
  Association Perspective</atl>
<src>Hospital and health Services
  Administration</src>
<dt>May&ndash;June 1982</dt>
</ft>
</tbl>
    
```

The above tagged table describes contents and structure. The variety of presentations for printing must be catered for with either a T_EX format or a L^AT_EX style file.

(Procedural) T_EX markup, AAP’s simple table

We have limited ourselves to “translating” SGML markup into procedural T_EX markup (no L^AT_EX markup). (Mainly: <name> into \bname and </name> into \ename ; the tags can be transformed table driven, but in the header row \nl commands have to be inserted manually, guided by taste and aesthetics and limited by the value of \entrywidth. Note that the data have to be adapted too: insertion of ? for suppressed 0, and \before %).

```

\entrywidth=2cm
\tablewidth=4\entrywidth
\btbl %AAP’s simple example with TeX
  %procedural markup
\bno Table AAP\eno
\bt Job Changes: 1973--1980 \ett
\bhd
  \bth\eth
  \bth Net Gain/Loss\nl
    of Hospitals \nl
    since 1973 \eth
  \bth Total No. \nl
    of CEO \nl
    Job Changes \nl
    1973--80 \eth
  \bth Survival \nl
    Rate of \nl
    CEO’s \eth
\ehd
\btby
\brow\btsb Texas\etsb\bc$+$20\ec
  \bc\ec\bc 22\%\ec
\erow
\brow\btsb Maryland\etsb\bc$+$?5\ec
  \bc 42\ec\bc 24\%\ec
\erow
\etby
\bsrc
  \bau David Kinzer\eau
  \batl Turnover Of Hospital Chief
  Executive Officers:
  A Hospital Association Perspective\eatl
  \bnme Hospital and health Services
  Administration\enme
  \bdt May--June 1982\edt
\esrc
\etbl
    
```

Note. We still have to supply the values for entrywidth and tablewidth along with each particular table, once again, manually.

T_EX format macros

In order to reproduce AAP's representation the following (format) macros were written

```
%TeX ``format'' for AAP's simple table.
\newdimen\entrywidth%\entrywidth=<default>
\newdimen\tablewidth
\tablewidth=\hsize%default
\def\nl{\par\noindent}
\def\ndash{--}
%? command for nonsignificant
% leading zeroes, Knuth p241
\catcode'?'=\active
\def?\{\kern\digitwidth}
%
\def\tablerule{\noalign{\hrule}}
\def\btbl{\bgroup
  \def\bno##1\eno{{\bf##1}}
  \def\btt##1\ett{{##1}\hidewidth\cr}
  \def\bhd{\tablerule\noalign{\vskiplex}}
  \def\ehd{\cr
    \noalign{\vskip.5ex}\tablerule
    \noalign{\vskip.5ex}}
  \def\bth##1\eth{\vtop{\noindent
    \hsize=\entrywidth ##1}&}
  \def\btby{\noalign{\vskiplex}}
  \def\etby{\noalign{\vskiplex}}
  \def\btsb##1\etsb{\hbox to
    \entrywidth{##1\hss}\hfil}
  \def\bc##1\ec{\&\hbox to .5
    \entrywidth{\hss ##1}\hfil}
  \def\brow##1\erow{##1\cr}
  \def\bsrc{\noalign\bgroup}
  \def\esrc{%Source information is
    %handled conform AAP's
    %representation
    Source: \gau, ``\gat1,``
    {\it \gnme\}
    \gdt.\ \gobi
    \egroup}
  % Next items are ``stored'' via gdefs
  \def\bau##1\eau{\gdef\gau{##1}}
  \def\bat1##1\eat1{\gdef\gat1{##1}}
  \def\bnme##1\enme{\gdef\gnme{##1}}
  \def\bdt##1\edt{\gdef\gdt{##1}}
  $$\vbox\bgroup\hsize=\tablewidth
  \halign\bgroup &##\cr%preamble line
  \tablerule\noalign{\vskiplex}
}%end\btbl
\def\etbl{\egroup%\halign
  \egroup$$\vbox
  \egroup%\btbl
}%end \etbl
%end AAP simple table format
```

The above listed format macros take care of the final results in print: appropriate separators and good order and format of the 'source' items. The table entries are centered and aligned on the last digit. This required

knowledge of the column width. In order to handle the footer suitably the tablewidth had to be known. The chosen approach allows flexible formatting of the footer. Variability of column widths has not been incorporated in the provided macros but can be addressed.

Difficulties with AAP's complex table DTD

Although not the case in the above elaborated example, we experienced difficulties with header rows which contain "halfines." In my opinion, halfines belong to the structure. Confusion arises when the br (begin row) and er (end row) attributes are used together with halfines. According to the DTD, halfines don't account for a line in the formatted result, in T_EX formatting they do, jeopardizing the prescribed br- and er-values.

Note that author etc. information is stored in gdefs in T_EX, in order to cope with the difference in sequence order of this items in SGML (AAP's DTD) and T_EX (independent) marked up tables.

Graphics

Coupling graphics is not (yet) elaborated, because graphics in SGML is left to other sources. CGM (Computer Graphics Metafile) methodology is the way SGML would like to see graphics incorporated. (Comm. Malcolm Clark, Idle by the thames, T_EXline X.) Various graphic sources are interfaced to SGML. The only structuring aspect deals with open space (to (electronically) paste up the illustration) which must not be split over a page break. The difference with mathematics possibly is that formulas are also part of the text while illustrations are more or less separated from the text.

Developments

A survey of the development of SGML is given by Barron [1990].

Usage

Some uses of SGML would be:

- DOD (Automated Technical Order System)
- European Communities (FORMalised EXchange of Electronic Documents; office official publications)
- Publishers (AAP, British Library, KNUB(Elsevier, Kluwer, ...), ...)
- Her Majesty's Stationary Office (legal text)
- HP Technical documentation
- Oxford University Press (abridged forms, database applications)
- McGraw Hill Encyclopedia of Science and technology (CD-ROM)
- SGML Users Group (chapters in various countries)
- T_EX to SGML assisting tools: XTRAN, Texttagger, Fasttag (communicated by Eric van Herwijnen)

Plans

Some plans for SGML use are being formulated:

- DOD (Computer-aided Acquisition and Logistic Support)
Object: To produce an integrated system in which information is held electronically, and which interfaces to CAD/CAM systems, electronic publishing systems and databases and those operated by the many defense contractors who supply the department, so that it will be possible to receive, distribute and use technical information in digital form.
- TEI-project, (Text Encoding and Interchange of machine readable texts), Sperberg-McQueen & Burnard [1990].

Local work in progress

Some local SGML work that is being done is:

- Elsevier's experiment, Bleeker [1989] and Poppelier [1990].
- Tabular material examples, van der Laan, Coleman [1990].
- Coupling SGML to L^AT_EX (mathematics) Grootenhuis [1990].

Guidelines for Choosing

As always whether to choose to use SGML, L^AT_EX, or T_EX depends upon many factors:

what is the document like,
what tools are already in use,
for whom is it aimed at,
how many authors are involved
is (partial, e.g. bibliographical) reuse also the case,
is future use, different from formatting, in sight?

No structure For documents with little structure, just standard font use, and no page make-up complexity, it does not matter what is used, provided minimal quality is obtained.

Scientific papers For scientific papers with complex mathematical, or physical, as well as tabular structures, T_EX/L^AT_EX can best be used, especially when publishers accept T_EX marked copy.

Reuse When reuse is the issue, e.g. bibliographic information, or document parts stored in a database, SGML can best be used.

Various authors When various authors with diverse backgrounds and text processing tools, provide copy to a publisher, it is tempting to consider SGML as a uniform language. Suitable tools

with proof facilities are necessary in order to stimulate authors to use it. One can also consider T_EX which is *de facto* in use for that purpose.

Future (nonformatting) use Abstract from medium, medium neutrality, and use SGML standard.

And remember: "SGML is on our minds, T_EX is in our hands."

What to do Next?

What about (digital) sound and (digital) video as part of the structured information?

Hypertext(, SGML(, T_EX(, . . .) . . .)?

While still struggling with DTDs, formatting, printing and coupling problems, it is tempting to ponder once in a while about *Information Type Definitions*.

Conclusion

SGML may benefit from T_EX formatting, and T_EX may benefit from SGML descriptive markup. From the mathematical and tabular examples it can be seen that there are difficulties in transforming SGML marked up copy into T_EX, when purism is strived for. The pragmatic approach with T_EX as a generally accepted SGML NOTATION, for (displayed) mathematical and tabular T_EX marked up copy, seems practical. Agreed this limits to formatting, but T_EX is stable and will serve a life-time. T_EX marked up copy can always be converted when needed. This approach has also been mentioned by TEI (Sperberg-McQueen & Burnard, 1990). Warmer mentions also difficulties in writing DTDs for tabular material because of the dichotomy of the tree structure: based on rows or on columns.

Acknowledgements

Most SGML codings are tentative, only the original SGML codings of mathematics have been parsed. Mary and Dean Guenther are kindly acknowledged for their support in molding the article into acceptable size and for changing the text into palatable English.

References

- [1] Appelt, W.(1988): Introduction to T_EX. Addison-Wesley.
- [2] Association of American Publishers (1987): Standard for electronic manuscript preparation and markup. AAP inc. ⁸
- [3] Association of American Publishers (1989): Author's guide to electronic manuscript preparation and markup. AAP inc. (2nd version.)

⁸ Association of American Publishers, 2005 Massachusetts Avenue, NW. Washington, DC 20036, Phone: (202) 232-3335

- [4] Association of American Publishers (1989): Reference manual on electronic manuscript preparation and markup. AAP inc. (2nd version.) ISBN 1-55653-084-6.
- [5] Association of American Publishers (1989): Markup of tabular material. AAP Inc. (2nd version.) ISBN 1-55653-085-4.
- [6] Association of American Publishers (1989): Markup of Mathematical Formulas. AAP Inc. EPSIG. (2nd version.)
- [7] Barron, D.(1989): Why use SGML? Electronic publishing, 2,1, 3–24.
- [8] Bechtolsheim, S. von (1990): T_EX in practice. (4 volumes, 1200p.).
- [9] Bleeker, J.(1989): Electronische verzending, bewerking en opslag van wetenschappelijke artikelen. In: SGML de consequenties. De eerste Nederlandse SGML conferentie. SGML User's Group Holland. (Dutch)
- [10] Bryan, M.(1988): SGML, an Author's Guide to the Standard Generalized Markup Language. Addison-Wesley.
- [11] Buerger, D.(1990): L^AT_EX for scientists and engineers, McGraw-Hill.
- [12] Cheswick, B.(1990): A permuted index for T_EX and L^AT_EX. CSR 145. AT&T Bell laboratories.
- [13] Clark, M.(1988): A note comparing T_EX to SGML. SGML User's Group Bulletin, 3, 2, 67–68.
- [14] Clark, M. (1989): T_EX and/or SGML. Proceedings EuroT_EX89. Karlsruhe. (Context sensitivity as a tool for checking input correctness is stressed; an example of how to do this within T_EX is given.)
- [15] Coombs, J.H., A.H. Renear, S.J. DeRose (1987): Markup systems and the future of scholarly text processing. Comm. ACM, 30, 11, 933–947.
- [16] Doob, M.(1989): A gentle introduction to T_EX. A manual for selfstudy.
- [17] Einarsson, B.(1988): Mixed Language programming. Sotware-Practice and Experience.
- [18] Genussa, P.L.(1987): Document Preparation Method of the United States Aire Force Automated Technical Order System. SGML Users' group. Bulletin 2, 1.
- [19] Grootenhuis, J.(1990): Typesetting SGML coded Mathematics. Paper presented at Markup'90. Charleston.
- [20] Guittet, C.(1986): FORMEX: une mise en pratique des normes internationales. SGML user's group. Bulletin, 1, 2.
- [21] Herwijnen, E. van (1988): Electronic submission of Physics articles to publishers. De 1^e Nederlandse SGML conferentie. SGML: De Consequenties. (Also published in:Computer Physics Communications 57 (1989) 244–250: The use of text interchange standards for submitting physics articles to journals.). In the context of this paper the discussion of SGML related to T_EX is relevant.)
- [22] Herwijnen, E. van (1990): Practical SGML. Kluwer-Academic.
- [23] ISO8879 Information Processing—Text and Office Systems—Standard Generalized Markup Language (SGML). 1986–10–15.
- [24] ISO/TR9573 Information Processing—SGML support facilities —Techniques for using SGML. 1988–09–12.
- [25] Knuth, D.E. (1989): The Errors of T_EX. Softw. prac. exp. 19, 7. 607–685.
- [26] Laan, C.G. van der (1984): (Graceful) Mixed Language Programming. Argonne National Laboratory Workshop.
- [27] Laan, C.G. van der (1990): Typesetting Bridge via T_EX. TUGboat, 11#2, 265–276.
- [28] Laan, C.G. van der, D.C. Coleman, J.R. Luyten (1989): SGML–(L^A)T_EX. 1. Mathematical Formulas. RC-RUG report 24.
- [29] Laan, C.G. van der, D.C. Coleman (to appear): SGML–(L^A)T_EX. 2. Tabular material.
- [30] L^AT_EX (1985): L^AT_EX a Document Preparation System. Addison-Wesley.
- [31] Poppelier, N.A.F.M.(1990): SGML and T_EX in Scientific Publishing. EuroT_EX90, Cork.
- [32] Seroul, R.(1989): Le petit livre de T_EX. Inter-Éditions. Paris.
- [33] MARK-IT (1989): SGML Parser, version 2. Sobemap NV, Place du Champ de Mars 5, Bte 40, 1050 Bruxelles.
- [34] Schwarz, N.(1989): Einführung in T_EX. Addison-Wesley. (Translated into Dutch and English)
- [35] Smith, J.M.(1987): The standard generalized markup language (SGML): Guidelines for editors and publishers. British National Bibliography Research Fund Report 26. ISBN 0-7123-3111-5.
- [36] Smith, J.M.(1987): The standard generalized markup language (SGML): Guidelines for authors. British National Bibliography Research Fund Report 27. ISBN 0-7123-3112-3.
- [37] SGML User's Group Newsletters. ⁹

⁹Editorial address: Pindar Infotek, 2 Grosvenor Road, Wallington, Surrey SM6 0ER, UK.

- [38] Sperberg-McQueen, C.M., L. Bernard (1990): ACH-ACL-ALLC. Guidelines for the encoding and interchange of machine readable texts. 2, 2, 65–90.
- [39] Vignaud, D.(1989): L'éditior structrée dans les documents, SGML applications à l'éditior française. Éditions du Cercle de la Librairie. Paris.
- [40] Warmer, J., S. van Egmond (1989): The implementation of the Amsterdam SGML Parser. EP-ODD, [41] Warmer, J. (priv. comm).
- [42] Wittbecker, A.(1989): T_EX enslaved. Proceedings TUG89. Stanford. (Advantages and disadvantages of T_EX-formatter with SGML “front-end” are discussed, related to DEC's VAX Document.)