

GELLMU: A Bridge for Authors from L^AT_EX to XML

Presentation at TUG 2001, University of Delaware

William F. Hammond

Department of Mathematics & Statistics
University at Albany
Albany, New York 12222 (USA)
hammond@math.albany.edu
<http://www.albany.edu/~hammond/>

August 2001
(slightly revised: January 2003)

Abstract

GELLMU, which stands for “Generalized Extensible L^AT_EX-Like Markup”, is a system for using L^AT_EX-like markup, though not L^AT_EX itself, to write consciously for a markup language in the SGML category or in its popular XML subcategory.

The *basic* level of GELLMU offers a way to use L^AT_EX-Like notation together with a L^AT_EX-Like *newcommand* (with arguments) macro facility to write web pages.

The *advanced* level of GELLMU enables one additionally to incorporate certain L^AT_EX-Like features, such as the use of a blank line for a new paragraph, in writing for an SGML language.

The didactic GELLMU production system provides an “article” XML language, with some resemblance to L^AT_EX itself, that is a rigorous domain for translation to other formats.

Table of Contents

| | | |
|---|---|---|
| 1 | Author Level Markup | 2 |
| 2 | The basics of <i>basic</i> GELLMU | 3 |
| 3 | Basic GELLMU enhanced with <code>\newcommand</code> | 4 |
| 4 | SGML vs. XML | 4 |
| 5 | Advanced GELLMU | 5 |
| 6 | The GELLMU didactic production system | 6 |
| 7 | Production of this Document | 7 |

1 Author Level Markup

Inasmuch as the World Wide Web is becoming an important library resource, one wants one's publications to be accessible online, and one wants web-crawling robots to be able to catalogue them properly.

Despite the popularity of Adobe's Portable Document Format (PDF) the distribution of PDF reading software is not as widespread as the distribution of web browsing software, and web-crawling robots often do not scan the contents of PDF documents.

What is available for the \LaTeX author toward this end? More specifically, consider the following situations:

Online publication archives Specifically, I would like to cite the \TeX / \LaTeX -based e-print archive begun at Los Alamos in the early 1990's by Paul Ginsparg, now known as "Arxiv"¹, which is now a participant in the Open Archives Initiative². While in its early time the term *e-print* was understood to mean "electronic pre-print", ArXiv has more recently become a repository for established journals including now, for example, the highly regarded *Annals of Mathematics*³, which was founded in 1884 by Ormond Stone of the University of Virginia, and Ginsparg now tells us that the term *e-print* denotes "self-archiving by the author" under a new overall academic publication⁴ design.

Course handouts How can a college teacher prepare course handouts for both paper and online distribution? If the teacher writes \LaTeX , some manual intervention will likely be needed in order to obtain correct HTML. If the teacher writes HTML, then the paper distribution⁵ will be limited by what can be expressed in HTML, which is not as rich a markup as \LaTeX .

TUG articles Before preparing a TUG 2001 article an author is asked to read *Preparation of documents for multiple modes of delivery*⁶ by Ross Moore, which is available on the web only as a two-column PDF printed page image. From this article one might conclude that carefully prepared \LaTeX may be suitable for translation to HTML although no HTML version of the article seems to be available.

GNU documentation While working for TUG on the \TeX Directory System (TDS) guidelines — see `/tds/standard/`⁷ at CTAN — in January 1998, Ulrik Vieth produced a \LaTeX document and a tailored program for translating that document into *Texinfo*, the language of the GNU⁸ documentation system.

Texinfo is a \TeX -based system that pre-dates HTML. Its original purpose was to provide both print and (early online hypertext) *Info* versions of GNU software project documentation. When HTML came along, it was possible to provide fairly reliable translation

¹URI: <http://www.arxiv.org/>

²URI: <http://www.openarchives.org/>

³URI: <http://www.math.princeton.edu/~annals/>

⁴Paul Ginsparg, "Electronic Clones vs. the Global Research Archive", <http://arXiv.org/blurb/pg00bmc.html>.

⁵If the HTML is correctly written, then robust translation to \LaTeX is possible.

⁶URI: <http://www.tug.org/TUGboat/Contents/contents20-4.html>

⁷URI: <http://ctan.tug.org/tex-archive/tds/standard/>

⁸URI: <http://www.gnu.org/>

from *Texinfo* to HTML because *Texinfo* is a well-structured markup. In fact, *Texinfo* is very nearly equivalent to an SGML language, and, Daniele Giacomini in August 2000 came up with an effort in that direction: *Sgmltexi*⁹.

Although programs are available for translating carefully structured L^AT_EX into HTML and sometimes into XML extensions of HTML, this method of generating online content for the basic level of the web sometimes requires manual intervention. A more direct approach to the world of SGML offers better prospects for long-term access to new web formats without sacrificing access to the quality of print typesetting that is available through L^AT_EX.

2 The basics of *basic* GELLMU

In looking over Vieth's set-up for the TDS document in the late spring of 1998, I arrived at the idea of using L^AT_EX-like notation for conscious writing in document languages under SGML and I have written a program in the GNU Emacs Lisp language, the GELLMU syntactic translator, for converting this L^AT_EX-like markup to SGML markup.

The advantage of SGML markup is that each markup language (formally document type) under the SGML umbrella constitutes a structured domain for the application of automatic processors that are easy to create under any of a number of structured processing frameworks. There are frameworks accessible in standard computing languages, and there is also a recent framework *xmltex*¹⁰ by David Carlisle for writing T_EX typesetting routines for XML document types.

The root idea in using L^AT_EX-like markup for the conscious writing of markup under SGML is the simple syntactic correspondence between markup such as

```
some \em{emphasized} text
```

on the one hand, and the markup

```
some <em>emphasized</em> text
```

on the other.

Most L^AT_EX commands are analogous to SGML elements. Moreover, the attribute list associated with an SGML element can be made to correspond with a L^AT_EX command option. For example,

```
\a[href="http://foo.dom/"
]{The Foo Domain}
```

matches

```
<a href="http://foo.dom/"
>The Foo Domain</a>
```

⁹URI: <http://master.swlibero.org/~daniele/software/sgmltexi/>

¹⁰URI: <http://www.dcarlisle.demon.co.uk/xmltex/>

3 Basic GELLMU enhanced with `\newcommand`

The idea of using L^AT_EX-like syntax for conscious writing under an SGML document type gains power when one realizes that although the notion of SGML entity provides, among other things, simple macro expansions, there is no provision under SGML for macros that take arguments. Moreover, there is no obvious method of extending SGML systems to accommodate macros with arguments apart from the idea of extending a document type¹¹.

GELLMU provides a L^AT_EX-like meta-command¹² called `newcommand` that may be invoked with arguments. For example, if one writes

```
\newcommand{\afoo}[2]{%
  \a[href="http://www.foo.org/#1"]{#2}}
```

then a subsequent invocation

```
\afoo{tex-archive/tds/}{TDS at Foo}
```

will yield (without line breaks):¹³

```
<a href=
"http://www.foo.org/tex-archive/tds/"
>TDS at Foo</a>
```

This `newcommand` markup differs from that of L^AT_EX in that it is classical macro substitution rather than vocabulary expansion. Since the syntax of a `newcommand` invocation is very similar to that of an SGML element, the use of `newcommand` can, apart from its on-the-fly convenience, be a help in the development of SGML document type extensions. A new name in a test document can be moved from being that of a macro to being that of an element simply with the removal of a `newcommand` definition.

4 SGML vs. XML

Basic GELLMU as enhanced by its macro facility is as far as one can sensibly go toward conscious writing under a language in the restricted subfamily of SGML document types known as XML.

From one viewpoint the differences between SGML and XML are not very important since most correct documents under the larger category can, if correct, be automatically translated into equivalent documents under XML. For example, classical HTML that passes validation can be translated into the newer XML form of HTML using either James Clark's classical

¹¹While document type extensions require enough work that they cannot be spontaneous, they provide a sound way to avoid the tangles that can arise working with T_EX or L^AT_EX when attempting the simultaneous use of conflicting macro packages.

¹²In GELLMU while a `command` corresponds to an SGML element, a `meta-command` is something having the same syntax as a command that does not correspond to an SGML element and instead receives resolution into other SGML markup under the syntactic translator.

¹³The syntactic translator maintains line number alignment between its input and its SGML output so that line numbers used by SGML parsers in flagging errors match those in source markup.

SGML library *SP*¹⁴ or Dave Raggett's program *tidy*¹⁵.

However, the rules for XML were designed to make things easy for processors rather than for humans, and for that reason an author writing toward an ultimate XML document type usually is well-advised to write for a version of the document type under more author-friendly SGML rules.

For example, if in an SGML language a forced linebreak is represented by the defined-empty element *brk*, then the markup `<brk>` is sufficient, whereas under the more restrictive XML version of the same language, either the markup `<brk></brk>` or its abbreviated form `<brk/>`¹⁶ must be used. For GELLMU this means that the markup `\brk` and the markup `\brk{}` are interchangeable under SGML, except for the case of `\brk` abutting a following character without intervening whitespace, but not equivalent under XML. GELLMU provides the form `\brk;` to represent the abbreviated form `<brk/>` of an element that is defined as empty.

5 Advanced GELLMU

Basic GELLMU deals with markup languages more or less at the level of syntax without getting to the level of grammar.

Advanced GELLMU may be used to roll language-independent grammatical concepts into the picture.

The first of these is L^AT_EX-like multiple argument/option syntax. For example under advanced GELLMU the markup

```
\frac{a x + b}{c x + d}
```

is converted in syntactic translation to

```
<frac><ag0>a x + b</ag0><ag0>c x + d</ag0>
```

That is, a chain of '{', '}' pairs and '[', ']' pairs following a command without intervening white space between the command name and the first delimiter nor between a close delimiter and the next open delimiter in the chain, constitutes an SGML element whose content begins with a sequence of generic positional arguments (tag name *ag0*) and options (tag name *op0*). Without knowledge of the document type it cannot be determined if a name used with multiple argument/option syntax has only *ag0*, *op0* content. The syntactic translator provides a list variable consisting of names that have only this type of content and that, therefore should be given close tags after the sequence of arguments and options. Absent that, the author must provide a close tag unless an SGML parser can infer it, and even in that case, if the element can appear in the mixed content model of another element such as, for example, a paragraph, then the parser's automatic placement of a close tag could lead to the unwanted collapse of a word boundary similar to that which occurs in L^AT_EX when an author's careless markup

\TeX benchmark

¹⁴URI: <http://www.jclark.com/sp/>

¹⁵URI: <http://www.w3.org/People/Raggett/tidy/>

¹⁶Moreover, some confusion may arise from the fact that under the SGML syntax (formally SGML declaration) specified for HTML neither of these XML forms of markup would be permitted.

gets typeset as “ \TeX benchmark” instead of as “ \TeX benchmark”.

If multiple argument/option syntax is used, then there is ambiguity on the nature of the first pair of chained delimiters if it is the pair ‘ $[$ ’, ‘ $]$ ’ — whether it represents an *op θ* or an attribute list. Therefore, in this case it is required that it is an attribute list if the first character after its ‘ $[$ ’ opening delimiter is a colon (‘ $:$ ’).

In basic GELLMU the following four of the ten \LaTeX -special characters are special:

$\backslash \{ \} \% .$

Additionally, the character ‘ $\#$ ’ is special when used in the definition of a *newcommand*, the characters ‘ $[$ ’ and ‘ $]$ ’ are special when used for \LaTeX -like option syntax, and the character ‘ $\&$ ’ is special when followed immediately by a letter since then it is the introducer for SGML entity invocations.

Advanced GELLMU provides for the possibility of giving traditional \LaTeX meaning to ‘ $\&$ ’ when not followed immediately by a letter and also to the other four \LaTeX -special characters, which are

$_ \^ \$ \sim .$

Additionally, it provides for the possibility of giving traditional \LaTeX -like meaning to other short forms of markup such as “ $\backslash(. . . \backslash)$ ” for inline math, “ $\backslash[. . . \backslash]$ ” for displayed math, “ $--$ ” for a range-dash, “ $---$ ” for a punctuation-dash, “ \backslash ” for an inter-word space, “ $\backslash,$ ” for a short space, and others including also, if desired, the use of blank lines, as appropriate, for paragraph boundaries.

6 The GELLMU didactic production system

The conversion of both *basic* and *advanced* GELLMU source markup to SGML is performed by my program called the *syntactic translator*.

If one wishes to write consciously for a public document type such as HTML or the Text Encoding Initiative¹⁷’s TEI using GELLMU’s \LaTeX -like syntax, the syntactic translator is the only part of GELLMU that will be of interest.

The optional features of advanced GELLMU described above can only be used when one is writing for a document type that provides markup in which the corresponding concepts have representation. For example, \LaTeX -like use of the character ‘ \sim ’ for non-breaking space requires a markup that provides non-breaking space.

Moreover, if blank lines are going to be paragraph boundaries, then the syntactic translator will need a list of element names before which a new paragraph does not make sense, and, since there is no separate provision for a list of names after which a paragraph must end, the document type cannot be XML.

¹⁷URI: <http://www.tei-c.org>

The GELLMU didactic production system provides such a document type and also provides tools, which can be used as inter-changeable components, for working with that document type.

The didactic production system consists of the syntactic translator and the following additional components:

1. An SGML document type called “article”.
2. Its XML cousin, also called “article”.
3. A program for translating SGML article to XML article.
4. A program for translating XML article to HTML.
5. A program for translating XML article to \LaTeX .

The document type is intended to be comfortable for authors with past experience in \LaTeX . The document type and the components are didactic. They are intended to illustrate how such a system can be assembled from inter-changeable components. They are not finished in any sense, and each has shortcomings.

They do serve, I hope, to demonstrate to the community of \LaTeX authors that it will find no limitations in this approach to document production.

At the same time it is intended to provide a whole new way of thinking about the subjects of package design and class design.

Its unfinished nature is intended to make it relatively easy for those who are so inclined to move in various ways to finish such a system that fits their needs.

7 Production of this Document

This document and the slides used during its presentation were prepared with the GELLMU didactic production system. Pre-publication versions of the sources and various automatic formattings are available in the author’s web¹⁸.

Subsequent to the GELLMU run on this document a copy of its \LaTeX output was manually modified for conformance with TUG guidelines. If I were going to submit a number of such TUG articles, it would be worthwhile to make another variant of the \LaTeX formatter for TUG.

¹⁸URI: <http://math.albany.edu:8000/math/pers/hammond/Presen/tug2001/>