

TEX Users Group

Et φθΑεηπw82ωLPR9BΘκAΨΓΙλS
l Δ N V 4 Ω W 0 l e 5 0 B X h 3 6 X 0 3 π Y K M 7 L 2
Z C K L O O Ψ I E 6 0 A T m 7 σ Γ B a λ E V 3 n
Graphic artists, typographers, programmers and
TEXnicians, all of you with a child's heart,
come and join us at the
**International
Conference on
TEX, XML and
Digital
Typography
TUG2004**
in Xanthi, Greece,
between August 30
and September 3,
2004!

Hosted by the Greek TEX
Friends Group in Xanthi,
Greece, in collaboration
with the Democritus
University of Thrace.

For more information,
write to us at:
tug2004@ocean1.ee.duth.gr
or visit our Web site:
obelix.ee.duth.gr/tug2004

STATHIS

05AKWXEΞQ1n0εf x m d ⊕ τ ρ ∞ Ho

Drawing by Stathis Stavropoulos
Poster design by "TeX by Hand Unlimited"
© Copyright 2003 Greek TEX Friends

PREPRINTS
for the 2004 Annual Meeting

T_EX Users Group

These preprints for the 2004 annual meeting are published by the T_EX Users Group.

Periodical-class postage paid at Portland, OR, and additional mailing offices. Postmaster: Send address changes to T_EX Users Group, 1466 NW Naito Parkway Suite 3141, Portland, OR 97209-2820, U.S.A.

Memberships

2004 dues for individual members are as follows:

- Ordinary members: \$75.
- Students/Seniors: \$45.

The discounted rate of \$45 is also available to citizens of countries with modest economies, as detailed on our web site.

Membership in the T_EX Users Group is for the calendar year, and includes all issues of *TUGboat* for the year in which membership begins or is renewed, as well as software distributions and other benefits. Individual membership is open only to named individuals, and carries with it such rights and responsibilities as voting in TUG elections. For membership information, visit the TUG web site: <http://www.tug.org>.

Institutional Membership

Institutional Membership is a means of showing continuing interest in and support for both T_EX and the T_EX Users Group. For further information, contact the TUG office (office@tug.org) or see our web site.

T_EX is a trademark of the American Mathematical Society.

Copyright © 2004 T_EX Users Group.

Copyright to individual articles within this publication remains with their authors, and may not be reproduced, distributed or translated without their permission.

For the editorial and other material not ascribed to a particular author, permission is granted to make and distribute verbatim copies without royalty, in any medium, provided the copyright notice and this permission notice are preserved.

Permission is also granted to make, copy and distribute translations of such editorial material into another language, except that the T_EX Users Group must approve translations of this permission notice itself. Lacking such approval, the original English permission notice must be included.

Printed in U.S.A.

Board of Directors

Donald Knuth, *Grand Wizard of T_EX-arcana*[†]
Karl Berry, *President*^{*}
Kaja Christiansen^{*}, *Vice President*
Sam Rhoads^{*}, *Treasurer*
Susan DeMeritt^{*}, *Secretary*
Barbara Beeton
Jim Hefferon
Ross Moore
Arthur Ogawa
Gerree Pecht
Steve Peter
Cheryl Ponchin
Michael Sofka
Philip Taylor
Raymond Goucher, *Founding Executive Director*[†]
Hermann Zapf, *Wizard of Fonts*[†]

^{*}member of executive committee

[†]honorary

Addresses

General correspondence,
payments, etc.
T_EX Users Group
P. O. Box 2311
Portland, OR 97208-2311
U.S.A.
Delivery services,
parcels, visitors
T_EX Users Group
1466 NW Naito Parkway
Suite 3141
Portland, OR 97209-2820
U.S.A.

Telephone

+1 503 223-9994

Fax

+1 503 223-3960

Electronic Mail

(Internet)
General correspondence,
membership, subscriptions:
office@tug.org
Submissions to *TUGboat*,
letters to the Editor:
TUGboat@tug.org
Technical support for
T_EX users:
support@tug.org
Contact the Board
of Directors:
board@tug.org

World Wide Web

<http://www.tug.org/>

Problems not resolved?

The TUG Board wants to hear from you:
Please email board@tug.org.

[printing date: April 2004]

**PREPRINTS
for the 2004 Annual Meeting**

TeX Users Group
Twenty-fifth Annual Meeting
Xanthi, Greece
August 30 – September 3, 2004

A COMMUNICATION OF THE T_EX USERS GROUP

EDITORS KARL BERRY
 BADEN HUGHES
 STEVE PETER

PORTLAND • OREGON • U.S.A. • 2004

Introduction

Welcome to this special communication from TUG! Although it is formatted in the usual *TUGboat* proceedings style, as you can see, it is not officially a *TUGboat* issue, there is no volume or issue number, and it's being produced *before* the meeting. Why? These editorial notes (try to) explain the situation.

The 2004 TUG Annual Meeting (<http://tug.org/tug2004>) is being organized by Apostolos Syropoulos and others, associated with the Greek T_EX Friends user group, in collaboration with the Democritus University of Thrace in Xanthi, Greece. The organizers were able to negotiate with Springer-Verlag for the proceedings of the meeting to be published in Springer's prestigious series *Lecture Notes in Computer Science*.

As TUG president, I felt such a publication was in the best interests of T_EX and TUG, as it would expose our system to a wide new audience. On the other hand, it is of course absolutely necessary for every current TUG member to receive the proceedings, as that is a benefit promised upon joining.

Unfortunately, it is not financially possible for TUG to purchase the hardbound Springer volume for every member. Therefore, with the consent of the TUG board, we worked out this compromise. The principal points of the agreement are these:

- The present volume of "preprints" is fully edited and prepared to usual *TUGboat* publication standards.
- The present volume will be one of the usual three TUG publications sent to members for 2004. The other two will be regular *TUGboat* issues.
- Springer's volume will be the publication of record for the 2004 meeting; this is why the present volume is not an official *TUGboat* issue.
- We also are not sending the present volume to libraries, as it cannot easily be cataloged. Thus, it would only cause trouble for the librarians. Libraries will receive the Springer volume.

- TUG will not sell this issue through our online store (<http://tug.org/store>); it will be sent only to those who join TUG for the 2004 membership year.
- Springer or the authors may make the papers available online; for any that are available, we may link to them from the *TUGboat* web pages, but we won't publish them online ourselves.

I sincerely hope this one-time special arrangement will meet with everyone's understanding. We expect future proceedings of TUG annual meetings to be published as normal *TUGboat* issues.

The theme of the meeting is *T_EX in the era of Unicode*, and this is reflected in the contents here. There is a wealth of material on Omega and other multilingual typesetting, including a definitive study of the historical Deseret alphabet. Additional papers include a comparative study of Type 1 font generation from METAFONT sources, a significant new package for critical edition typesetting, and more.

As this volume is being prepared in advance of the meeting, please email any comments you might have directly to the authors, or to the editors at tug2004@ocean1.ee.duth.gr, so the presentation might be improved. And of course please consider going to the conference. (Online registration is available from the conference web site.)

Production of the three issues of *TUGboat* for the 2003 membership year is continuing. The proceedings of the TUG 2003 conference will probably be the first to be completed.

As always, feel free to email the TUG board (board@tug.org), or the *TUGboat* editorial staff (tugboat@tug.org) with any questions, concerns, or (most important of all) suggestions and submissions for future articles.

◇ Karl Berry
president@tug.org

Digital Typography in the New Millennium: Flexible Documents by a Flexible Engine

Christos KK Loverdos

Department of Informatics and Telecommunications

University of Athens

TYPA Buildings, Panepistimioupolis

GR-157 84 Athens

Greece

loverdos@di.uoa.gr

Apostolos Syropoulos

Greek T_EX Friends Group

366, 28th October Str.

GR-671 00 Xanthi

Greece

apostolo@obelix.ee.duth.gr

Abstract

The T_EX family of electronic typesetters is the primary typesetting tools for the preparation of demanding documents, and have been in use for many years. However, our era is characterized, among others, by Unicode, XML and the introduction of interactive documents. In addition, the Open Source movement, which is breaking new ground in the areas of project support and development, enables masses of programmers to work simultaneously. As a direct consequence, it is reasonable to demand the incorporation of certain facilities to a highly modular implementation of a T_EX-like system. Facilities such as the ability to extend the engine using common scripting languages (e.g., Perl, Python, Ruby, etc.) will help in reaching a greater level of overall architectural modularity. Obviously, in order to achieve such a goal, it is mandatory to attract a greater programming audience and leverage the Open Source programming community. We argue that the successful T_EX-successor should be built around a microkernel/exokernel architecture. Thus, services such as client-side scripting, font selection and use, output routines and the design and implementation of formats can be programmed as extension modules. In order to leverage the huge amount of existing code, and keep document source compatibility, the existing programming interface is demonstrated to be just another service/module.

1 Introduction

The first steps towards computer typesetting took place in the 1950s, but it was not until Donald E. Knuth introduced T_EX in 1978 [16] that true quality was brought to software-based typesetting. The history of T_EX is well-known and the interested reader is referred to [16] for more details.

Today, the original T_EX is a closed project in the sense that its creator has decided to freeze its development. As a direct consequence no other programs are allowed to be called T_EX. In addition, the freely available source code of the system was a major step on the road towards the formation of the

Open Source movement, which, in turn, borrowed ideas and practices from the Unix world. Furthermore, the development of T_EX and its companion system, METAFONT, had made obvious the need for properly documented programs. This, in turn, initiated Knuth's creation of the *literate programming* program development methodology. This methodology advances the idea that the program code and documentation should be intermixed and developed simultaneously.

The source code of T_EX and METAFONT being freely available has had enormous consequences. Anyone can not only inspect the source code, but

also experiment freely with it. Combined with \TeX 's (primitive, we should note, but quite effective for the time) ability to extend itself, this led to such success stories as \LaTeX and its enormous supporting codebase, in the form of packages. As a direct consequence of the fact that the source code is frozen, stability was brought forth. Note that this was exactly the intention Knuth had when developing his systems. A common referred-to core, unchanged in the passing of time and almost free of bugs, offered a "secure" environment to produce with and even experiment with.

However, in an everchanging world, especially in the fast-paced field of computer science, almost anything must eventually be surpassed. And it is the emerging needs of each era that dictate possible future directions. \TeX has undoubtedly served its purpose well. Its Turing-completeness has been a most powerful asset/weapon in the battles for and of evolution. Yet, the desired abstraction level, needed to cope with increasing complexity, has not been reached. Unfortunately, with \TeX being bound to a fixed core, *it cannot be reached*.

Furthermore, the now widely accepted user-unfriendliness of \TeX as a language poses another obstacle to \TeX 's evolution. It has created the myth of those few, very special and quite extraordinary "creatures"¹ able to decrypt and produce code fragments such as the following:²

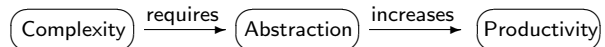
```
\def\s@vig{{\EO@m=\EO@n
\divide\EO@n by20 \relax
\ifnum\EO@n>0\s@vig\fi
\EO@k=\EO@n\relax
\multiply\EO@k by-20\relax
\advance\EO@m by \EO@k\relax
\global\advance\EO@l by \@ne
\expandafter\edef\csname EO@d\@roman{\EO@l}\endcsname{\fi
\ifnum\EO@m=0\noexpand\noexpand\EOzero
\else\expandafter\noexpand
\expandafter\csname EO\@roman{\EO@m}\endcsname\fi}
\expandafter\@rightappend
\csname EO@d\@roman{\EO@l}\endcsname
\t@{\epi@lme@Digits}}
```

Of course, to be fair, programmers in several languages (C and Perl among others) are often accused of producing understandable code and the well-known *obfuscated code* contests just prove it. On the other hand, with the advent of quite sophisticated assemblers, today one can even write well-structured assembly language, adhering even to "advanced" techniques/paradigms, such as object-oriented programming. Naturally, this should not lead to the conclusion that we should start writing in assembly (again)! In our opinion, software com-

¹ The second author may be regarded as one of Gandalf's famuli, while the first author is just a Hobbit, wishing to have been an Elf.

² Taken from the documentation of the *epi@lme@* package by the second author.

plexity should be tackled with an emphasis on abstraction that will eventually lead to increased productivity, as is shown in the following figure:



\TeX 's programming language is more or less an "assembly language" for electronic typesetting. It is true that higher level constructs can be made—macros and macro packages built on top of that. But the essence remains the same. Although it is true that \TeX is essentially bug free and its macro expansion facility behaves the way it is specified (i.e., as defined in [9]), it still remains a fact that it takes a non-specialist quite some time to fully understand the macro expansion rules in spite of Knuth's initial intentions [12, page 6].

The fact that one should program in the language of his/her choice is just another reason for moving away from a low-level language. And it is true that we envision an environment where as many programmers as possible can—and the most important, wish to—contribute. In the era of the Open Source revolution, we would like to attract the Open Source community and not just a few dedicated low-level developers. Open Source should also mean, in our opinion, "open possibilities" to evolve the source. This is one of our major motivations for reengineering the most successful typesetting engine.

Richard Palais, the founding chairman of TUG, pointed out back in 1992 [12, page 7] that when developing \TeX , Knuth

... had NSF grant support that not only provided him with the time and equipment he needed, but also supported a team of devoted and brilliant graduate students who did an enormous amount of work helping design and write the large quantity of ancillary software needed to make the \TeX system work ...

and immediately after this, he poses the fundamental question:

Where will the resources come from what will have to be at least an equally massive effort? And will the provider of those resources be willing, at the end of the project, to put the fruits of all his effort in the Public Domain?

The answer seems obvious now. The way has been paved by the GNU/Linux/BSD revolutionary development model, as has been explained crystal clearly in *The Cathedral and the Bazaar* [15].

This paper is an attempt to define a service-oriented architecture for a future typesetting engine, which will be capable of modular evolution.

We take a layered approach of designing some core functionality and then define extensible services on top of the core. The engine is not restricted to a specific programming language either for its basic/bootstrapping implementation or, even more important, for its future enhancement. At the same time, we are bound to provide a 100% \TeX -compatible environment, as the only means of supporting the vast quantity of existing \TeX -based documents. We intend to achieve such a goal by leveraging the proposed architecture's own flexibility. Specifically, a \TeX compatibility mode is to be supported and it should give complete "trip-test" compliance. Later on, we shall see that this compatibility is divided into two parts: source code compatibility and internal core compatibility. Both are provided by plug-gable modules.

Structure of the paper In the following sections we briefly review the most important and influential approaches to extending or reengineering \TeX , including \TeX 's inherent abilities to evolve. Then we discuss a few desired characteristics for any next generation typesetting engine. We advance by proposing an architecture to support these emerging needs. Finally, we conclude by discussing further and future work.

2 A Better \TeX ?

2.1 \TeX the Program

\TeX supports a Turing-complete programming language. Simply, this means that if it lacks a feature, it can be programmed. It contains only a few concepts and belongs to the LISP family of languages. In particular, it is a list-based macro-language with late binding [5, Sec. 3.3]:

Its data constructs are simpler than in Common Lisp: 'token list' is the only first order type. Glue, boxes, numbers, etc., are engine concepts; instances of them are described by token lists. Its lexical analysis is simpler than CL: One cannot program it. One can only configure it. Its control constructs are simpler than in CL: Only macros, no functions. And the macros are only simple ones, one can't compute in them.

Further analysis of \TeX 's notions and inner workings such as *category codes*, \TeX 's *mouth* and *stomach* is beyond the scope of this paper and the interested reader is referred to the classic [9] or the excellent [3].

\TeX the program is written in the WEB system of literate programming. Thus, its source code is self-documented. The programs `tangle` and `weave` are used to extract the Pascal code and the documentation, respectively, from the WEB code. The

documentation is of course specified in the \TeX notation. Although the \TeX source is structured in a monolithic style, its architecture provides for some kind of future evolution.

First, \TeX can be "extended" by the construction of large collections of macros that are simply called *formats*. Each format can be transformed to a quickly loadable binary form, which can be thought of as a primitive form of the module concept.

Also, by the prescient inclusion of the `\special` primitive command, \TeX provides the means to express things beyond its built-in "comprehension". For example, \TeX knows absolutely nothing about PostScript graphics, yet by using `\special` and with the appropriate driver program (e.g., `dvips`), PostScript graphics can be easily incorporated into documents. Color is handled in the same way. In all cases, all that \TeX does is to expand the `\special` command arguments and transfer the command to its normal output, that is, the DVI file (a file format that contains only page description commands).

Last, but not least, there is the notion of *change file* [3, page 243]:

A change file is a list of changes to be made to the WEB file; a bit like a stream editor script. These changes can comprise both adaptations of the WEB file to the particular Pascal compiler that will be used and bug fixes to \TeX . Thus the `TeX.web` file needs never to be edited.

Thus, change files provide a form of incremental modification. This is similar to the `patch` mechanism of Unix.

Yet, no matter how foresighted these methods may be, twenty years after its conception \TeX has started to show its age. Today's trends, and more importantly the programming community's continuing demand for even more flexible techniques and systems, call for new modes of expressiveness.

2.2 The \LaTeX Format

\LaTeX [10], which was released around 1985, is the most widely known \TeX format. Nowadays, it seems that \LaTeX is the de facto standard for the communication and publication of scientific documents (i.e., documents that contain a lot of mathematical notation). \LaTeX "programs" have a Pascal-like structure and the basic functionality is augmented with the incorporation of independently developed collections of macro packages. In addition, classes are used to define major document characteristics and are in essence document types, such as *book*, *article*, etc. Thus, each \LaTeX "program" is characterized by the document class to which it belongs, by the

packages it utilizes, and any new macro commands it may provide.

The current version of \LaTeX is called $\LaTeX 2_\epsilon$. Work is in progress to produce and widely distribute the next major version, $\LaTeX 3$ [11]. Among the several enhancements that the new system will bring forth, are:

- Overall robustness
- Extensibility, relating to the package interface
- Better specification and inclusion of graphical material
- Better layout specification and handling
- Inclusion of requirements of hypertext systems

The $\LaTeX 3$ core team expects that a major reimplementaion of \LaTeX is needed in order to support the above goals.

The $\text{Con}\TeX t$ [13] format, developed by Hans Hagen, is monolithic when compared to \LaTeX . As a result, the lessons learned from its development are not of great interest to our study.

2.3 $\mathcal{N}\mathcal{T}\mathcal{S}$: The New Typesetting System

The $\mathcal{N}\mathcal{T}\mathcal{S}$ project [14] was established in 1992 as an attempt to extend \TeX 's typesetting capabilities and at the same time to propose a new underlying programmatic model. Its originators recognised that \TeX lacked user-friendliness and as a consequence it attracted many fewer users than it could (or should). Moreover, \TeX (both as a name and a program) was frozen by Knuth, so any enhancements should be implemented in a completely new system.

$\mathcal{N}\mathcal{T}\mathcal{S}$ was the first attempt to recognize that \TeX 's monolithic structure and implementation in an obsolete language (i.e., the Pascal programming language) are characteristics that could only impede its evolution. The techniques used to implement \TeX , particularly its "tight", static and memory conservative data structures have no (good) reason to exist today (or even when $\mathcal{N}\mathcal{T}\mathcal{S}$ was conceived, in 1992), when we have had a paradigm shift to flexible programming techniques.

After considering and evaluating several programming paradigms [19] including functional, procedural and logic programming, the $\mathcal{N}\mathcal{T}\mathcal{S}$ project team decided to proceed with a Java-based implementation. Java's object-oriented features and its network awareness were the main reasons for adopting Java, as $\mathcal{N}\mathcal{T}\mathcal{S}$ was envisioned as a network-based program, able to download and combine elements from the network.

Today, there is a Java codebase, which has deconstructed the several functional pieces of \TeX and reconstructed them in a more object-oriented way

with cleaner interfaces, a property that the original \TeX source clearly lacks. In spite of the promising nature of $\mathcal{N}\mathcal{T}\mathcal{S}$, the directory listing at CTAN³ shows that the project is inactive since 2001.⁴ It seems that the main focus is now the development of $\epsilon\text{-}\TeX$, which is presented in the following section.

2.4 $\epsilon\text{-}\TeX$

$\epsilon\text{-}\TeX$ [17] was released by the $\mathcal{N}\mathcal{T}\mathcal{S}$ team as soon as it was recognized that $\mathcal{N}\mathcal{T}\mathcal{S}$ itself was very ambitious and that a more immediate and more easily conceivable goal should be set. So, it was decided that the first step towards a new typesetting system was to start with a reimplemented but 100% \TeX compatible program.

$\epsilon\text{-}\TeX$ was released in 1996, after three years of development and testing. It adds about thirty new primitives to the standard \TeX core, including handling of bidirectional text (right-to-left typesetting). It can operate in three distinct modes:

1. "compatibility" mode, where it behaves exactly like standard \TeX .
2. "extended" mode, where its new primitives are enabled. Full compatibility with \TeX is not actually sought and the primary concern is to make typesetting easier through its new primitives.
3. "enhanced" mode, where bidirectional text is also supported. This mode is taken to be a radical departure from standard \TeX .

Today, $\epsilon\text{-}\TeX$ is part of all widely used \TeX distributions and has proven to be very stable. Indeed, in 2003 the \LaTeX team requested that future distributions use $\epsilon\text{-}\TeX$ by default for \LaTeX commands, which has since been implemented in \TeX Live and other distributions.

2.5 Ω

Ω [16], which was first released in 1996, is primarily the work of two people: Yannis Haralambous and John Plaice. It extends \TeX in order to support the typesetting of multilingual documents. Ω provides new primitives and new facilities for this reason. Ω 's default character encoding is the Unicode UCS-2 encoding, while it can easily process files in almost any imaginable character encoding. In addition to that, Ω supports the parameterization of paragraph and page direction, thus allowing the typesetting of text in any imaginable writing method.⁵

³ <http://www.ctan.org/tex-archive/systems/nts/>.

⁴ We have last accessed the above URL in March 2004.

⁵ Currently the the boustrophedon writing method is the only one not supported.

Much of its power comes from its new notion of Ω TPs (Ω Translation Processes). In general, an Ω TP is normally used to transform a document from a particular character encoding to another. Obviously, an Ω TP can be used to transform text from one character set to another. An Ω TP is actually a finite state automaton and, thus, it can easily handle cases where the typesetting of particular characters are context dependent. For example, in traditional Greek typography, there are two forms of the small letter theta, which are supported by Unicode [namely ϑ (03D1) and θ (03B8)]. The first form is used at the beginning of a word, while the second in the middle of a word. The following code borrowed from [16] implements exactly this feature:

```
input: 2; output: 2;
aliases:
LETTER = (@"03AC-@"03D1 | @"03D5 | @"03D6 |
          @"03F0-@"03F3 | @"1F00-@"1FFF) ;
expressions:
~({LETTER})@"03B8({LETTER} | @"0027)
=> \1 @"03D1 \3;
. => \1;
```

For performance reasons, Ω TPs are compiled into Ω CPs (Ω Compiled Processes).

External Ω TPs are programs in any programming language that can handle problems that cannot be handled by ordinary Ω TPs. For example, one can prepare a Perl script that can insert spaces in a Thai language document. Technically, external Ω TPs are programs that read from the standard input and write to the standard output. Thus, Ω is forking a new process to allow the use of an external Ω TP. In [16] there are a number of examples (some of them were borrowed from [7]).

We should note that the field of multilingual typesetting is an active research field, which is the main reason why Ω is still an experimental system. We should also note that ε - Ω [4], by Giuseppe Bilotta, is an extension of Ω that tries to incorporate the best features of ε - $\text{T}_{\text{E}}\text{X}$ and Ω in a new typesetting engine.

2.6 pdf $\text{T}_{\text{E}}\text{X}$

pdf $\text{T}_{\text{E}}\text{X}$ [18] is yet another $\text{T}_{\text{E}}\text{X}$ extension that can directly produce a file in Adobe's PDF format. Recently, pdf- ε - $\text{T}_{\text{E}}\text{X}$ was introduced, merging the capabilities of both pdf $\text{T}_{\text{E}}\text{X}$ and ε - $\text{T}_{\text{E}}\text{X}$.

3 Towards a Universal Typesetting Engine

From the discussion above, it is obvious that there is a trend to create new typesetting engines that provide the best features of different existing typesetting engines. Therefore, a Universal Typesetting

Engine should incorporate all the novelties that the various $\text{T}_{\text{E}}\text{X}$ -like derivatives have presented so far. In addition, such a system should be designed by taking into serious consideration all aspects of modern software development and maintenance. However, our departure should not be too radical, in order to be able to use the existing codebase. Let us now examine all these issues in turn.

3.1 Discussion of Features

Data Structures $\text{T}_{\text{E}}\text{X}$'s inherent limitations are due to the fact that it was developed in a time when computer resources were quite scarce. In addition, $\text{T}_{\text{E}}\text{X}$ was developed using the now outdated structured programming program development methodology.

Nowadays, hardware imposes virtually no limits in design and development of software. Also, new programming paradigms (e.g., aspect-oriented programming [8], generative programming [2], etc.) and techniques (e.g., extreme programming [1]) have emerged, which have substantially changed the way software is designed and developed.

These remarks suggest that a new typesetting engine should be free of "artificial" limitations. Naturally, this is not enough as we have to leave behind the outdated programming techniques and make use of modern techniques to ensure the future of the Universal Typesetting Engine. Certainly, $\mathcal{N}\mathcal{T}\mathcal{S}$ was a step in the right direction, but in the light of current developments in the area of software engineering it is now a rather outdated piece of software.

New Primitive Commands Modern document manipulation demands new capabilities that could not have been foreseen at the time $\text{T}_{\text{E}}\text{X}$ was created. A modern typesetting engine should provide a number of new primitive commands to meet the new challenges imposed by modern document preparation. Although the new primitives introduced by ε - $\text{T}_{\text{E}}\text{X}$ and Ω solve certain problems (e.g., bidirectional or, more generally, multidirectional typesetting), they are still unable to tackle other issues, such as the inclusion of audio and/or animation.

Input Formats For reasons of compatibility, the current input format must be supported. At the same time the proliferation of XML and its applications makes it more than mandatory to provide support for XML content. Currently, XML $\text{T}_{\text{E}}\text{X}$ is a $\text{T}_{\text{E}}\text{X}$ format that can be used to typeset validated XML

files.⁶ In addition, $X\LaTeX$ [6] is an effort to reconcile the \TeX world with the XML world. In particular, $X\LaTeX$ is an XML Document Type Definition (DTD) designed to provide an XMLized syntax for \LaTeX . However, we should learn from the mistakes of the past and make the system quite adaptable. This means that as new document formats emerge, the system should be easily reconfigurable to “comprehend” these new formats.

Output Formats The $pdf\LaTeX$ variant has become quite widespread, due to its ability to directly produce output in a very popular document format (namely Adobe’s Portable Document Format). Commercial versions of \TeX are capable of directly generating PostScript files without the need of any driver programs. However, as in the case of the input formats, it is quite possible that new document formats will appear. Thus, we need to make sure that these document formats will find their way into \TeX sooner or later.

In addition, XML initiatives such as MathML and SVG (Scalable Vector Graphics) are increasingly common in electronic publishing of scientific documents (i.e., quite demanding documents from a typographical point of view). Thus, it is absolutely necessary to be able to choose the output format(s) from a reasonable list of options. For example, when one makes a drawing using \LaTeX ’s `picture` environment, it would be quite useful to have SVG output in addition to the “standard” output. Currently, Ω can produce XML content, but it cannot generate PDF files.

Innovative Ideas The assorted typesetting engines that follow \TeX ’s spirit are not mere extensions of \TeX . They have introduced a number of useful features and/or capabilities. For example, Ω ’s Ω T P s and its ability to handle Unicode input by default should certainly make their way into a new typesetting engine. In addition, ε - \TeX ’s new conditional primitives are quite useful in macro programming.

Typesetting Algorithms The paragraph breaking and hyphenation algorithms in \TeX make the difference when it comes to typographic quality. Robust and adaptable as they are, these algorithms may still not produce satisfactory results for all possible cases. Thus, it is obvious that we need a mech-

anism that will adapt the algorithms so they can successfully handle such difficult cases.

Fonts Typesetting means to put type (i.e., font glyphs) on paper. Currently, only `METAFONT` fonts and PostScript Type 1 fonts can be used with all different \TeX derivatives. Although Ω is Unicode aware, still it cannot handle TrueType fonts in a satisfactory degree (one has to resort to programs like `ttf2t1m` in order to make use of these fonts). In addition, for new font formats such as OpenType and SVG fonts there is only experimental support, or none at all. A new typesetting engine should provide font support in the form of plug-ins so that support for new font formats could be easily provided.

Scripting Scripting is widely accepted as a means of producing a larger software product from smaller components by “gluing” them together. It plays a significant role in producing flexible and open systems. Its realization is made through the so-called “scripting languages”, which usually are different from the language used to implement the individual software components.

One could advance the idea that scripting in \TeX is possible by using \TeX the language itself. This is true to some extent, since \TeX works in a form of “interpretive mode” where expressions can be created and evaluated dynamically at runtime—a feature providing the desired flexibility of scripting languages. But \TeX itself is a closed system, in that almost everything needs to be programmed within \TeX itself. This clearly does not lead to the desired openness.

A next generation typesetting engine should be made of components that can be “glued” together using any popular scripting language. To be able to program in one’s language of choice is a highly wanted feature. In fact, we believe it is the only way to attract as many contributors as possible.

Development Method Those software engineering techniques which have proven successful in the development of real-world applications should form the core of the program methodology which will be eventually used for the design and implementation of a next generation typesetting engine. Obviously, generic programming and extreme programming as well as aspect-oriented programming should be closely examined in order to devise a suitable development method.

All the features mentioned above as well as the desired ones are summarized in Table 1.

⁶ Validation should be handled by an external utility. After all, there are a number of excellent tools that can accomplish this task and thus it is too demanding to ask for the incorporation of this feature in a typesetting engine.

	\TeX	$\mathcal{N}\mathcal{T}\mathcal{S}$	$\varepsilon\text{-}\text{\TeX}$	Ω	$\text{\LaTeX}(3)$	Desired
implementation language	traditional	Java	traditional	traditional	traditional	perhaps scripting
architecture	monolithic	modular?	monolithic	monolithic	monolithic	modular
\TeX compatibility	100%	yes	100%	100%	100%	via module
input transformations				Ω TPs		pluggable
Unicode	(Babel)	(Java)	(Babel)	true		true
XML				yes	via package	yes
typesetting algorithms	\TeX	\TeX -like	\TeX -like	\TeX -like	\TeX -like	pluggable
scripting language	\TeX	$\mathcal{N}\mathcal{T}\mathcal{S}$ (?)	$\varepsilon\text{-}\text{\TeX}$	Ω	\TeX	any
output drivers	dvi(ps,pdf)	dvi(?)	dvi(ps,pdf)	dvi(ps,pdf)	dvi(ps,pdf)	any
TRIP-compatible	yes	almost	$\varepsilon\text{-}$ TRIP	yes	yes	yes (via module)
library mode	no	no	no	no	no	yes
daemon (server) mode	no	no	no	no	no	yes
programming community	< \LaTeX	1 person?	< \TeX	very small	big	> \LaTeX

 Table 1: Summary of features of \TeX and its extensions.

3.2 Architectural Abstractions

Roughly speaking, the *Universal Typesetting Engine* we are proposing in this paper, is a project to design and, later, to implement a new system that will support all the “good features” incorporated in various \TeX derivatives plus some novel ideas, which have not found their way in any existing \TeX derivative.

Obviously, it is not enough to just propose the general features the new system should have—we need to lay down the concrete design principles that will govern the development of the system. A reasonable way to accomplish this task is to identify the various concepts that are involved. These concepts will make up the upper abstraction layer. By following a top-down analysis, eventually, we will be in position to have a complete picture of what is needed in order to proceed with the design of the system.

The next step in the design process is to choose a particular system architecture. \TeX and its derivatives are definitely monolithic systems. Other commonly used system architectures include the microkernel and exokernel architectures, both well-known from operating system research.

Microkernel Architecture A microkernel-based design has a number of advantages. First, it is potentially more reliable than a conventional monolithic architecture, as it allows for moving the major part of system functionality to other components, which make use of the microkernel. Second, a microkernel implements a flexible set of primitives, providing high level of abstraction, while imposing little or no limitations on system architecture. Therefore, building a system on top of an existing microkernel is significantly easier than developing it from scratch.

Exokernel Architecture Exokernels follow a radically different approach. As with microkernels, they take as much out of the kernel as possible, but rather than placing that code into external programs (mostly user-space servers) as microkernels do, they place it into shared libraries that can be directly linked into application code. Exokernels are extremely small, since they arbitrarily limit their functionality to the protection and multiplexing of resources.

Both approaches have their pros and cons. We believe that a mixed approach is the best solution. For example, we can have libraries capable of handling the various font formats (e.g., Type 1, TrueType, OpenType, etc.) that will be utilized by external programs that implement various aspects of the typesetting process (e.g., generation of PostScript or PDF files). Let us now elaborate on the architecture we are proposing. The underlying components are given in Figure 1.

The *Typesetting Kernel* (TK) is one of the two core components at the first layer. It can be viewed as a “stripped-down” version of \TeX , meaning that its role as a piece of software is the orchestration of several typesetting activities. A number of basic algorithms are included in this kernel both as abstract notions—necessary for a general-purpose typesetting engine—and concrete implementations. So, TK incorporates the notions of paragraph and page breaking, mathematical typesetting and is Unicode-aware (utilizing UCS-4 internally). It must be emphasized that TK “knows” the concept of paragraph breaking and the role it plays in typesetting but it is not bound to a specific paragraph breaking algorithm. The same principle applies to all needed algorithms.

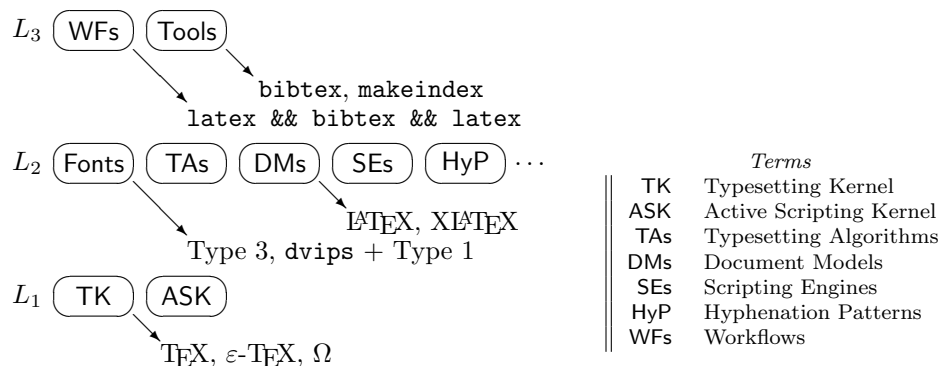


Figure 1: The proposed microkernel-based layered architecture. The arrows show rough correspondence between the several architectural abstractions and their counterparts in existing monolithic typesetting engines.

The *Active Scripting Kernel* (ASK) is the second of the core components and the one that allows scripting at various levels, using a programming (scripting) language of one’s choice. It is in essence a standardized way of communicating between several languages (T_EX, Perl, Python), achieved by providing a consistent Application Programming Interface (API). The most interesting property of ASK is its *activeness*. This simply means that any extension programmed in some language is visible to any other available languages, as long as they adhere to the standard Active Scripting Kernel API. For example, an external module/service written in Perl that provides a new page breaking algorithm is not only visible but also available for immediate use from Python, C, etc.

Above TK and ASK, at the second layer, we find a collection of typesetting abstractions.

Fonts are at the heart of any typesetting engine. It is evident that font architectures change with the passing of time, and the only way to allow for flexibility in this part is to be open. Although there are many different font formats, all are used to define glyphs and their properties. So instead of directly supporting all possible font formats, we propose the use of an abstract font format (much like all font editors have their own internal font format). With the use of external libraries that provide access to popular font formats (e.g., a Free Type library, a Type 1 font library, etc.), it should be straightforward to support any existing or future font format.

The various *Typesetting Algorithms* (TAs)—algorithms that implement a particular typographic feature—should be coded using the Active Scripting Kernel API. In a system providing the high degree of flexibility we are proposing, it will be possible to exhibit, *in the same document*, the result of applying

several paragraph and page breaking algorithms. By simply changing a few runtime parameters it will be possible to produce different typographic “flavors” of the same document.

A *Scripting Engine* (SE) is the realization of the ASK APIs for a particular scripting language. For reasons of uniformity, the T_EX programming language will be provided as a Scripting Engine, along with engines for Perl, Ruby and Python. This will make all the existing T_EX codebase available for immediate use and it will provide for cooperation between existing L^AT_EX packages and future enhancements in other languages. Thus, a level of 100% T_EX compatibility will be achieved, merely as a “side-effect” of the provided flexibility.

The idea of a *Document Model* (DM) concerns two specific points: The document *external* representation, as it is “edited” for example in an editor, or “saved” on a hard disk, and its *internal* representation, used by the typesetting engine itself. It is clear that under this distinction, current L^AT_EX documents follow the (fictional) “L^AT_EX Document Model”, X_LL^AT_EX documents follow the “X_LL^AT_EX document model” and an XML document with its corresponding DTD follows an analogous “XML+DTD Document Model”.

We strongly believe that how a document is written should be separated by its processing. For the last part, an *internal* representation like the *Abstract Syntax Trees* (ASTs) used in compiler technology is highly beneficial. One way to think of DM is as the typographic equivalent of the Document Object Model (DOM). That is, it will be a platform-neutral and language-neutral representation allowing scripts to dynamically access and update the content, structure and style of documents.

Several *Document Processors* (DPs) may be applied to a specific document before actual typesetting takes place. DPs are the analog of Ω TPs. By leveraging the scripting power of ASK, the representation expressiveness of DPs is increased—as opposed to algorithmic expressiveness (Turing-completeness), which is evident, e.g., in Ω , but is not the sole issue.

The *Workflows* (WF) and *Tools* are at the highest architectural layer. Currently, there are a number of tools that may not produce a final typeset result, but are important for the proper preparation of a document. For example, such tools include bibliography, index and glossary generation tools. In the proposed architecture, all these programs will take advantage of other architectural abstractions—such as the Document Model or the Scripting Engines—in order to be more closely integrated in the typesetting engine as a whole.

Of particular importance is the introduction of the *Workflows* notion. A workflow is closely related to the operation or, to be more precise, *cooperation* of several tools and the typesetting engine in the course of producing a typeset document. In effect, a workflow specifies the series of execution (probably conditional) steps and the respective inputs/outputs during the “preparation” of a document. By introducing a workflow specification for each tool, we relieve the user from manually specifying all the necessary actions in order to get a “final” .pdf (or whatever output format has been requested). Instead, the user will declaratively specify that the services of a tool are needed and the engine will load the respective workflows, compose them and execute them.

We shall give a workflow example concerning a `bibtex`-like tool. What we do here is to transform our experience of using `bibtex` into declarations specifying its behaviour in cooperation with `latex`:

```
WORKFLOW DEFINITION bibtex

SERVICE bibtex NEEDS latex
SERVICE bibtex INTRODUCES latex
```

In effect, this translates a hypothetical `Makefile`:

```
all:
    latex mydoc
    bibtex mydoc
    latex mydoc
```

for the preparation of the fictitious `mydoc.tex` document into a declarative specification that is given only *once* as part of the `bibtex` tool!

3.3 On Design and Evolution

Recent advances in software engineering advocate the use of multidimensional separation of concerns as a guiding design principle. Different concerns should be handled at different parts of code and ideally should be separated. For example, the representation of a document and its processing are two separate concerns and should be treated as such. Their interaction is better specified out of their individual specifications. Thus, we have introduced the Document Models notion to cope with the existing \TeX / \LaTeX base as well as any future document representation.

Several architectural abstractions of Figure 1 are candidates to be specified as “services” at different granularities. For example, any *Tool* of the third layer can be thought of as a service that is registered with a naming authority and discovered dynamically, for immediate use on demand. A True-Type Font Service, regarding the second layer *Font* abstraction, is another example, this time more of a fine-grained nature, in the sense that a *Tool* (coarse-grained service) utilizes a *Font* (fine-grained service).

The proposed architecture makes special provisions for evolution by keeping rigid design decisions to a minimum. Built-in Unicode awareness is such a notable rigid design decision, but we feel that its incorporation is mandatory. Besides that, the ideas of pluggable algorithms and scripting are ubiquitous and help maintain the desired high degree of flexibility.

At the programming level, any style of design and development that promotes evolution can be applied. In the previous section we have actually demonstrated that the proposed architecture can even handle unanticipated evolution at the workflow level: the `bibtex` tool workflow specification causes the execution of an existing tool (`latex`) but we have neither altered any workflow for `latex` nor does `latex` need to know that “something new” is using it. In effect, we have *introduced* (the use of the keyword `INTRODUCE` was deliberate) a new *aspect* [8].

4 Conclusions and Future Work

In this paper we have reviewed the most widespread modern approaches to extending \TeX , THE typesetting engine. After analyzing weaknesses of the approaches and the existing support for several features, we have presented our views on the architecture of an open and flexible typesetting engine.

We have laid down the basic architectural abstractions and discussed their need and purpose. Of

course, the work is still at the beginning stages and we are now working on refining the ideas and evaluating design and implementation approaches.

The introduction of the Active Scripting Kernel is of prime importance and there is ongoing work to completely specify a) the form of a standard procedural API and b) support for other programming styles, including object-oriented and functional programming. This way, an *object* may for example take advantage of an algorithm that is better described in a *functional* form. There are parallel plans for transforming T_EX into a Scripting Engine and at the same time providing Engines powered by Perl and Python.

We are also investigating the application of the workflow approach at several parts in the architecture other than the interaction among tools. This, in turn, may raise the need for the incorporation of a *Workflow Kernel* at the core layer, along with the Typesetting Kernel and the Active Scripting Kernel.

References

- [1] chromatic. *Extreme Programming Pocket Guide*. O'Reilly & Associates, Sebastopol, CA, USA, 2003.
- [2] Krzysztof Czarnecki and Ulrich Eisenecker. *Generative Programming: Methods, Tools, and Applications*. Addison-Wesley Publ. Co., Reading, MA, USA, 2002.
- [3] Victor Eijkhout. T_EX by Topic. <http://www.cs.utk.edu/~eijkhout/tbt>.
- [4] ε - Ω Project home page. <http://www.ctan.org/tex-archive/systems/eomega/>.
- [5] $\mathcal{N}\mathcal{T}\mathcal{S}$ FAQ. <http://www.ctan.org/tex-archive/info/NTS-FAQ>.
- [6] Yannis Haralambous and John Plaice. Omega, OpenType and the XML World. *The 24th Annual Meeting and Conference of the TeX Users Group, TUG 2003*.
- [7] Yannis Haralambous and John Plaice. Traitement automatique des langues et compositions sous omega. *Cahiers GUTenberg*, pages 139–166, 2001.
- [8] Gregor Kiczales, John Lamping, Anurag Mendhekar, Chris Maeda, Cristina Lopes, Jean-Marc Loingtier, and John Irwin. Aspect-Oriented Programming. In M. Aksit and S. Matsuoka, editors, *ECOOP '97 — Object-Oriented Programming: 11th European Conference, Jyväskylä, Finland, June 1997. Proceedings*, number 1241 in Lecture Notes in Computer Science, pages 220–242. Springer-Verlag, Berlin, 1997.
- [9] Donald Erwin Knuth. *The T_EXbook*. Addison-Wesley, 1984.
- [10] Leslie Lamport. *L^AT_EX: A Document Preparation System*. Addison-Wesley Publ. Co., Reading, MA, USA, 2nd edition, 1994.
- [11] L^AT_EX3 Project home page. <http://www.latex-project.org/latex3.html>.
- [12] Richard Palais. Position Paper on the future of T_EX. <http://www.loria.fr/services/tex/moteurs/nts-9207.dvi>, reached from <http://tex.loria.fr/english/moteurs.html>, October 1992.
- [13] PRAGMA Advanced Document Engineering. ConT_EXt home page. <http://www.pragma-ade.com/>.
- [14] $\mathcal{N}\mathcal{T}\mathcal{S}$ Project home page. <http://www.dante.de/projects/nts/>.
- [15] Eric E. Raymond. The Cathedral and the Bazaar. <http://www.catb.org/~esr/writings/cathedral-bazaar/>.
- [16] Apostolos Syropoulos, Antonis Tsolomitis, and Nick Sofroniou. *Digital Typography Using L^AT_EX*. Springer-Verlag, New York, NY, USA, 2003.
- [17] $\mathcal{N}\mathcal{T}\mathcal{S}$ Team and Peter Breitenlohner. The ε -T_EX manual, Version 2. *MAPS*, (20):248–263, 1998.
- [18] Hàn Thé Thành, Sebastian Rahtz, and Hans Hagen. The pdfT_EX users manual. *MAPS*, (22):94–114, 1999.
- [19] Jiří Zlatuska. $\mathcal{N}\mathcal{T}\mathcal{S}$: Programming Languages and Paradigms. EuroT_EX 1999, <http://www.uni-giessen.de/partosch/eurotex99/zlatuska.pdf>.

Migrating to XML: The Case of the GUST Bulletin Archive

Włodzimierz Bzyl

Instytut Matematyki, Uniwersytet Gdański
80-952 Gdańsk
ul. Wita Stwosza 57
Poland
matwb@univ.gda.pl

Tomasz Przechlewski

Wydział Zarządzania, Uniwersytet Gdański
81-824 Sopot
ul. Armii Krajowej 119/121
Poland
tomasz@gnu.univ.gda.pl

Abstract

Ten years of experience with $\text{T}_{\text{E}}\text{X}$ publishing of the GUST bulletin shows that Knuth's dream of highly portable $\text{T}_{\text{E}}\text{X}$ files is apparently an illusion in practice. Over the last decade, articles in the GUST bulletin have used at least six major formats (\LaTeX 2.09, transitional \LaTeX +NFSS, \LaTeX 2 ϵ , plain-based *TUGboat*, Eplain, and Con $\text{T}_{\text{E}}\text{X}$ t), numerous macro packages, fonts, and graphic formats. Many old articles are typeset differently nowadays, and some even cause $\text{T}_{\text{E}}\text{X}$ errors.

This situation motivates the following question: how do we avoid the same problem in the future? As the World Wide Web is quickly becoming the mainstream both of publishing and of information exchange we argue for migration to XML—a Web compatible data format.

In the paper we examine a possible strategy for storing GUST articles in a custom XML format and publishing them with both $\text{T}_{\text{E}}\text{X}$ and XSLT/FO. Finally, the problems of converting the $\text{T}_{\text{E}}\text{X}$ files to XML and possibility of using $\text{T}_{\text{E}}\text{X}$ 4ht—an authoring system for producing hypertext—are discussed.

1 Introduction

The dominant role played by the Web in information exchange in modern times has motivated publishers to make printed documents widely available on the Internet. It is now common that many publications are available on the Web only, or before they are printed on paper. Articles published in the GUST bulletin are available on the Web in PostScript and PDF. Unfortunately, these formats decrease document accessibility, searching and indexing by Web search engines. For broad accessibility to automated services, it is better to use XML as the format of such data. However, one issue with XML is that it is difficult to maintain the high quality presentation of $\text{T}_{\text{E}}\text{X}$ documents. This is caused by incompatibilities between browsers and incomplete or immature implementations of W3C Consortium standards.

We are optimistic that these issues will disappear in the near future, and believe that XML will become pervasive in the online environment. However, in our context, a key to the adoption of XML is the degree to which it can be integrated with existing $\text{T}_{\text{E}}\text{X}$ nologies.

In this paper we examine one strategy for storing GUST articles in a custom XML format and publishing them with *both* $\text{T}_{\text{E}}\text{X}$ and XSLT/FO. Also, the problems of converting the existing $\text{T}_{\text{E}}\text{X}$ files to XML and the possibility of using $\text{T}_{\text{E}}\text{X}$ 4ht—an authoring system for producing hypertext—are discussed.

2 $\text{T}_{\text{E}}\text{X}/\text{\LaTeX}$ and Other Document Formats

When the authors started work with $\text{T}_{\text{E}}\text{X}$ (many years ago), there was only a choice between closed-source applications based on proprietary formats, or

\TeX , for publishing significant documents. Nowadays, the choice is much wider, as XML-based solutions are based on open standards and supported by a huge number of free applications. We do not need to write the tools ourselves. Thus the strategy of reusing what is publicly available is key in our migration plan.

On the other hand it would be unwise to switch to XML as the only acceptable submission format, because it would force many authors to abandon their powerful \TeX -based editing environments to which they are accustomed, just to submit texts to our bulletin. Following this route, we would more likely end up with a shortage of submissions. Thus, we are preparing a mixed strategy with both \TeX and XML as accepted formats. Papers submitted in \LaTeX will ultimately be converted to XML as an archival or retrieval format. Presentation formats will be XHTML, with corresponding PDF generated by a variety of tools. The work-flow of documents in this proposed framework is depicted on Fig. 1.

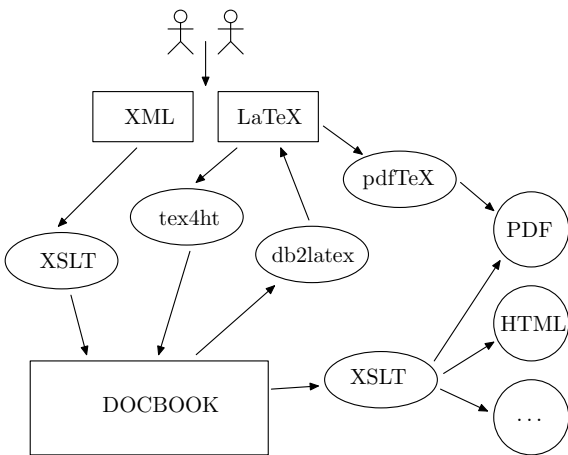


Figure 1: Processing diagram for XML/ \LaTeX documents.

The XML implementation project described in the paper can be broadly divided into the following subtasks: DTD development, formatting development, and legacy information conversion [19]. We'll now describe these stages in detail.

3 DTD Development Considerations

There is no doubt (see for example [14, 19]) that the DTD development phase is of critical importance for the overall success of any SGML/XML project. Fortunately, thanks to the great interest in XML technology in recent years, there are several production-quality publicly available DTDs which could be adapted for our project. To make this

choice, we preferred those which are widely used and for which formatting applications and tools are available. The following possible schemes were considered:

- DocBook [21], a complete publishing framework, i.e., schemes plus XSLT/DSSSL stylesheets for conversion to presentation formats; actively developed and maintained; the de facto standard of many Open Source projects; widely known and used.
- TEI [1], another complete, actively developed publishing framework. Not as popular as DocBook, used mainly in academia.
- The \LaTeX -based DTD developed in [7] (further referred as LWC DTD). The similarity to the structure of \LaTeX is an advantage of this DTD for our project.
- Others, such as DTD for GCA/Extreme conferences, X-DiMi from the Electronic Thesis and Dissertations Project, and the \LaTeX -like PMLA developed by one of the authors [15].

Industry-standard DTDs tend to be too big, too complex, and too general for practical use in specific cases (cf. [14, p. 29]). In particular, the DocBook and TEI DTDs seem to be too complex for marking-up documents conforming to \LaTeX format.

As a result, users frequently employ the technique of using different DTDs at different stages of the editorial process. Following Maler [14], we will call the DTD common to a group of users within an interest group as a *reference* DTD, while those used solely for editing purposes as an *authoring* DTD. Translation from one DTD to another may be easily performed with an XSLT stylesheet.

We decided to use a simplified LWC DTD as authoring DTD and DocBook as reference DTD. Providing a simple DTD should expand the group of prospective authors. For example, many GUST members are experts in typography or Web design but not necessarily \TeX hackers.

The simplification consists of restricting the document hierarchy only to article-like documents, and removing back matter tags (`index`, `glossary`) and all presentation tags (`newline`, `hspace`, etc.). Also, the optional status of meta-data, for example the `title`, `abstract`, `keywords` tags, was changed to required. The resulting DTD contains 45 elements compared to 64 in the original one.

For better maintainability, we rewrote our version of LWC DTD into RNC syntax. The RNC schema was introduced by Clark [6], and recently adopted as an ISO standard. It has many advantages

over DTD or W3C Schema syntax, namely simplicity and an included documentation facility.¹

As the structure of our documents is not particularly complex, it may be feasible to develop several authoring DTDs targeted at different groups of authors, for example one for those preferring ConTeXt-like documents, another for those used to GCA conference markup, etc., and then map those documents to the reference DTD with XSLT.

4 Formatting with XSLT

For presentation, LWC documents are first transformed to DocBook with a simple XSLT stylesheet.

The DocBook XSL stylesheets [22] translate an XML document to XHTML or FO [18]. As they are written in a modular fashion, they can be easily customized and localized. To publish XHTML from XML documents, an XSLT engine is needed such as Kay’s `saxon` [11] or Veillard’s `xsltproc` [20].

For hard copy output, a two-step process is used. First, the XSLT engine produces formatting objects (FO) which then must be processed with a formatting object processor for PDF output.² The detailed transformation work-flow is depicted in Fig. 2.

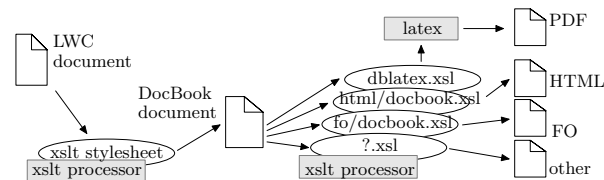


Figure 2: Processing details of LWC documents with XSLT/FO.

With just a few customizations the translation from XML to XHTML presents no obstacles (except for math formulas). On the other hand, the quality of the PDF produced with the publicly available `fop` processor from the Apache project is poor compared to that obtained with \TeX .

Instead of generating FO objects one can use XSLT to translate XML directly to high-level \LaTeX . That is the method used in `db2latex` [3] (see also a clone project: `dblatex/dbcontext` [9]); the latter, of course, generates files processable by ConTeXt).

¹ It is possible to convert between different schema languages for XML with the `trang` program [5]. There is also a `xml-mode` for GNU Emacs for editing of XML which features highlighting, indentation, and on the fly validation against an RNC schema [4].

² Modern browsers have XSLT engines built-in. So, it suffices to attach to a document appropriate stylesheets to make the transformation on the fly.

The output can be customized at XSLT stylesheet level as well as by redefining appropriate \LaTeX style files. MathML markup is translated with XSLT to \LaTeX and supported natively.³

The translation from DocBook to \LaTeX implemented in these tools is incomplete. To get reasonable results, prior customization to local needs is required. The main advantage of this approach is that we use \TeX —a robust and well known application.

5 The GUST Bulletin Archive

When considering the conversion of the GUST archive to XML we have two points in mind: first, we recognize the long-term benefits of an electronic archive of uniformly and generically marked-up documents; and second, to take the opportunity to test the whole framework using ‘real’ data.

During the last 10 years, 20 volumes of the GUST bulletin were published, containing more than 200 papers. From the very beginning GUST was tagged in a modified *TUGboat* style [2]. The total output is not particularly impressive, but the conversion of all articles to XML isn’t a simple one-night job for a bored \TeX hacker:

- they were produced over an entire decade and were written by over 100 different authors.
- they were processed with at least six major formats (\LaTeX 2.09, transitional \LaTeX +NFSS, \LaTeX 2 ϵ , plain-based *TUGboat*, Eplain, and finally ConTeXt), using numerous macro packages, fonts, and graphic formats.⁴

As a group, the GUST authors are not amateurs, producing naïve \TeX code. On the contrary they are \TeX experts, writing on a diverse range of subjects using non-standard fonts, packages and macros. For example, Fig. 3 shows the detailed distribution of the \TeX formats used in GUST.

In total, there were 134 plain \TeX articles, compared to 87 for \LaTeX . \LaTeX authors used 74 different packages, while those preferring plain \TeX nology used 139 different style files. The proportion of other formats (Eplain, ConTeXt, BLUE) was insignificant (only a few submissions). It can also be noted from Fig. 3 that in recent years, the proportion of plain \TeX submissions has decreased substantially in favor of \LaTeX .

It is obviously very difficult to maintain a repository containing papers requiring such a diverse

³ One approach which we did not try is to format FO files with \TeX . This method is implemented by S. Rahtz’ Passive \TeX [17].

⁴ Needless to say, all of these packages have been evolving during the last 10 years, many of them becoming incompatible with each other.

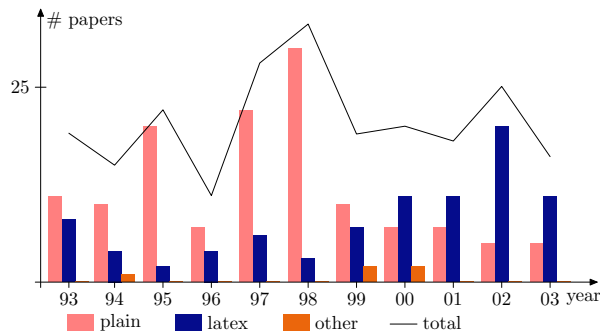


Figure 3: Distribution of TeX formats used by GUST authors.

range of external resources (macros, fonts). As a result, many old papers are now unable to be typeset owing to changes in underlying macros or fonts.

6 Conversion from TeX to XML

It may be surprising that only few papers report successful conversion from TeX to XML: Grim [8] describes successful large-scale conversion in a large academic institution, while Rahtz [16] and Key [12] describe translation to SGML at Elsevier.

Basically when converting TeX to XML the following three approaches have been adopted [16]:

- Perl/Python hackery combined with manual re-typing and/or manual XML marking-up.
- Parsing TeX source files not with `tex`, but with another program which generates SGML/XML. This is the approach used by `ltx2x` [23], `tralics` [8] and `latex2html`,⁵ which replace L^AT_EX commands in a document by user-defined strings.
- Processing files with TeX and post-processing the DVI files to produce XML. This is the way `tex4ht` works.

Although the conversion performed with `tralics` is impressive, we found the application very poorly documented. After evaluation of the available tools and consulting the literature [7], we decided to use TeX4ht—a TeX-based authoring system for producing hypertext [10].

Because TeX formats contain many visually-oriented tags, we could not expect to automatically convert them to content-oriented XML markup.⁶

For example, the *TUGboat* format requires only the metadata elements title and author name(s); author address(es) and webpage(s) of the

⁵ `latex2html` was not considered as its output is limited to HTML.

⁶ For example, see [16, 8]. Other examples, based on GUST articles, are presented below.

author(s) are often absent and there is no obligation for abstracts and keywords. Therefore, most of the GUST articles lack these valuable elements. Moreover, bibliographies are inconsistently encoded.⁷

Having said that, our plan is to markup as many elements as possible.

7 Translating TeX to XML with TeX4ht

Out of the box, the TeX4ht system is configured to translate from plain, L^AT_EX, *TUGboat* (`ltugboat`, `ltugproc`), and Lecture Notes in Computer Science (`llncs`) formats to HTML, XHTML, DocBook, or TEI. To translate from, say, *TUGboat* to our custom XML format the system needs to be manually configured. Because the configuration of TeX4ht from scratch is a non-trivial task, we consider other more efficient possibilities.

The TeX4ht system consists of three parts: (1) Style files which enhance existing macros with features for the output format (HTML, XHTML, etc.).⁸ (2) The `tex4ht` processor which extracts HTML or XHTML, or DocBook, or TEI files from DVI files produced by `tex`. (3) The `t4ht` processor which is responsible for translating DVI code fragments which need to be converted to pictures; for this task the processor uses tools available on the current platform.

As mentioned above, the conversion from a visual format to an information-oriented one cannot be done automatically. Let's illustrate this statement with the following example marked with plain TeX macros:⁹

```
\noindent {\bf exercise, left as an}
{\it adj\} {\ss Tech} Used to complete
a proof when one doesn't mind a
{\bf handwave}, or to avoid one entirely.
The complete phrase is: {\it The proof
\rm(or {\it the rest\}\rm) \it is left as an
exercise for the reader.\} This comment
has occasionally been attached to unsolved
research problems by authors possessed of
either an evil sense of humor or a vast
faith in the capabilities of their
audiences.\hangindent=1em
```

After translation of this fragment to XHTML by `tex4ht` we obtain:

```
<p class="noindent"><span
class="cmbx-10">exercise, left as an
</span><span class="cmti-10">adj
```

⁷ Publicly available tools (see [13] for example) can automatically mark up manually keyed bibliographies.

⁸ Altogether over 2.5M lines of TeX code. Compare this with 1M code of the L^AT_EX base macros.

⁹ The text comes from “The Project Gutenberg Etext of The Jargon File”, Version 4.0.0.

```

</span><span class="cmss-10">Tech
</span>Used to complete a proof when one
doesn't mind a <span
class="cmbx-10">handwave</span>, or to
avoid one entirely. The complete phrase
is: <span class="cmti-10">The proof
</span>(or <span class="cmti-10">the
rest</span>) <span class="cmti-10">is left
as an exercise for the reader.
</span>This comment has occasionally been
attached to unsolved research problems by
authors possessed of either an evil sense
of humor or a vast faith in the capabilities
of their audiences.</p>

```

and this could be rendered by a browser as:

exercise, left as an *adj* Tech Used to complete a proof when one doesn't mind a **handwave**, or to avoid one entirely. The complete phrase is: *The proof (or the rest) is left as an exercise for the reader.* This comment has occasionally been attached to unsolved research problems by authors possessed of either an evil sense of humor or a vast faith in the capabilities of their audiences.

We can see that `tex4ht` uses 'span' elements to mark up font changes. These visual tags could be easily remapped to logical ones unless fragments of text with different meaning are marked with the same tag. Here the tag `cmti-10` was used to tag both the short form 'adj' and the example phrase (shown in the green italic font). To tag them differently we need different `TeX` macros specially configured for `TeX4ht`. Note that the `\hangindent=1em` was ignored by `tex4ht`. This command could not be simulated, because hanging indentation is not supported by browsers.

So, the markup produced by the `tex4ht` program is not logical markup. To get logical markup the GUST format should be reworked and reconfigured for `TeX4ht`.

Instead of configuring `TeX4ht` we could use an XSLT stylesheet to remap elements referencing XML format. This could be an easier route than configuring the system from scratch, while some `TeX4ht` configuration could also help. So, a combination of the two methods is envisaged to provide the best results.

8 Conclusion and Future Work

We have not completed the conversion yet. However, based on the experience gained so far we can estimate that almost 70% of the whole archive should be converted with little manual intervention. Semi-automatic conversion of another 15% (34 papers) is possible, with prior extensive changes in markup. Conversion of remaining 15% is impossible or useless, where 'impossible' means the paper is easier

to retype than try to recompile and adjust `tex4ht` just for a particular single case, and 'useless' applies to papers demonstrating complicated graphical layouts, or advanced typesetting capabilities of `TeX`.

Although our system needs improvement—conversion of math is the most important remaining item to investigate—we are determined to start to use it in a production environment.

Finally, we note that many of our conclusions and methods are also applicable to `TUGboat`, because the format used for typesetting GUST bulletin differs only slightly from the one used for `TUGboat`.

References

- [1] Lou Burnard and C. M. Sperberg-McQueen. TEI lite: An introduction to text encoding for interchange. <http://www.tei-c.org/Lite/>, 2002.
- [2] Włodek Bzyl and Tomasz Przechlewski. An application of literate programming: creating a format for the Bulletin of the Polish TUG. *TUGboat*, 14(3):296–299, October 1993.
- [3] Ramon Casellas and James Devenish. DB2LaTeX XSL stylesheets. <http://db2latex.sourceforge.net>, 2004.
- [4] James Clark. NXML mode for the GNU Emacs editor. <http://www.thaiopensource.com/download>, 2003.
- [5] James Clark. Trang—multi-format schema converter based on RELAX NG. <http://www.thaiopensource.com/relaxng/trang.html>, 2003.
- [6] James Clark and Makoto Murata. Relax NG specification. <http://www.relaxng.org/>, 2001.
- [7] Michel Goossens and Sebastian Rahtz. *L^AT_EX Web Companion*. Addison-Wesley, 2001.
- [8] Jose Grim. Tralics. In *EuroTeX Preprints*, pages 38–49, 2003. <http://www-sop.inria.fr/miaou/tralics>. Final version to appear in *TUGboat*.
- [9] Benoît Guillon. DocBook to L^AT_EX/ConT_EXt publishing. <http://dblatex.sourceforge.net>, 2004.
- [10] Eitan Gurari. `tex4ht`: L^AT_EX and T_EX for hypertext. <http://www.cis.ohio-state.edu/~gurari/TeX4ht/mn.html>, 2004.
- [11] Michael Kay. SAXON—the XSLT and XQuery Processor. <http://saxon.sourceforge.net>, 2003.
- [12] Martin Key. Theory into practice: working with SGML, PDF and L^AT_EX. *Baskerville*,

- 5(2), 1995. ftp://tug.ctan.org/pub/tex/usergrps/uktug/baskervi/5_2/.
- [13] Language Technology Group. LT TTT version 1.0. <http://www.ltg.ed.ac.uk/software/ttt>, 1999.
- [14] Eve Maler and Jeanne El Andaloussi. *Developing SGML DTDs: From Text to Model to Markup*. Prentice Hall PTR, 1995.
- [15] Tomasz Przechlewski. Definicja dokumentu typu PMLA. <http://gnu.univ.gda.pl/~tomasz/sgml/pmla/>, 2002.
- [16] Sebastian Rahtz. Another look at L^AT_EX to SGML conversion. *TUGboat*, 16(3):315–324, September 1995. <http://www.tug.org/TUGboat/Articles/tb16-3/tb48raht.pdf>.
- [17] Sebastian Rahtz. PassiveT_EX. <http://www.tei-c.org.uk/Software/passivetex>, 2003.
- [18] Robert Stayton. *Using the DocBook XSL Stylesheets*. Sagehill Enterprises, <http://www.sagehill.net/docbookxsl/index.html>, 2003.
- [19] Brian E. Travis and Dale C. Waldt. *The SGML Implementation Guide: A Blueprint for SGML Migration*. Springer-Verlag, 1996.
- [20] Daniel Veillard. LIBXSLT—the XSLT C library for Gnome. <http://xmlsoft.org/XSLT>, 2003.
- [21] Norman Walsh and Leonard Muelner. *DocBook: The Definitive Guide*. O’Reilly, 1999. <http://www.docbook.org/tdg/en/html/docbook.html>.
- [22] Wiki. DocBook XSL Stylesheets. <http://docbook.org/wiki/moin.cgi/DocBookXslStylesheets>, 2004.
- [23] Peter R. Wilson. LTX2X: A L^AT_EX to X Auto-tagger. <http://www.ctan.org/tex-archive/support/ltx2x>, 1999.

Managing \TeX Resources with XML Topic Maps

Tomasz Przechlewski

Uniwersytet Gdański, Wydział Zarządzania

81-824 Sopot

ul. Armii Krajowej 119/121

Poland

tomasz@gnu.univ.gda.pl

Abstract

For many years the Polish \TeX Users Group newsletter has been published online on the GUST web site. The repository now contains valuable information on \TeX , METAFONT, electronic documents, computer graphics and related subjects. However, access to the content is very poor: it is available as PS/PDF files with only a simple HTML page facilitating navigation. There is no integration with information resources from other sections of the site, nor with the resources from other LUG or CTAN sites.

Topic maps were initially developed for efficient preparation of indices, glossaries and thesauruses for electronic documents repositories, and are now codified as both the ISO standard (ISO/IEC 13250) and the XTM 1.0 standard. Their applications extend to the domain of electronic publishing. Topic maps and the similar RDF standard are considered to be the backbone of corporate knowledge management systems and/or the Semantic Web [3].

The paper contains an introduction to the Topic Maps standard and discusses selected problems of Topic Map construction. Finally the application of Topic Maps as an interface to the repository of \TeX related resources is presented, as well as the successes and challenges encountered in the implementation.

1 Introduction

All the papers published for the last 10 years in the bulletin of the Polish \TeX Users' Group (GUST, <http://www.gust.org.pl/>) are now available online from the GUST web site. The repository contains valuable information on \TeX , METAFONT, electronic documents, computer graphics, typography and related subjects. However, access to the content itself is very poor: the papers are available as PS/PDF files with only a simple HTML interface facilitating navigation. There is no integration with other resources from that site. As CTAN and other LUGs' sites provide more resources it would obviously be valuable to integrate them too.

At first glance, the Topic Maps framework appears to be an attractive way to integrate vast amounts of dispersed \TeX related resources. A primary goal of the proposed interface should be to support learning. If the project succeeds, we hope it will change slightly the opinion of \TeX as a very difficult subject to become acquainted with.

The paper is organized as follows. The standard is introduced and selected problems of topic maps construction are discussed in the subsequent

three sections. Then a short comparison of Topic Maps and RDF is presented. The application of Topic Maps as an interface to the GUST resource repository is described in the last two sections.

2 What is a Topic Map?

Topic Maps are an SGML/HYTIME based ISO standard defined in [1] (ISO/IEC 13250, often referred as HYTM). The standard was recently rewritten by an independent consortium, TopicMaps.org [19] and renamed to XML Topic Maps (XTM). XTM was developed in order to simplify the ISO specification and enable its usage for the Web through XML syntax. Also, the original linking scheme was replaced by XLINK/XPOINTER syntax. XTM was recently incorporated as an Annex to [1].

The standard enumerates the following possible applications of TMs [1]:¹

- To qualify the content and/or data contained in information objects as topics, to enable navigational tools such as indexes, cross-references, citation systems, or glossaries.

¹ Examples of application of Topic Maps to real world problems can be found in [9, 21, 5, 18, 11, 12].

- To link topics together in such a way as to enable navigation between them.
- To filter an information set to create views adapted to specific users or purposes. For example, such filtering can aid in the management of multilingual documents, management of access modes depending on security criteria, delivery of partial views depending on user profiles and/or knowledge domains, etc.
- To add structure to unstructured information objects, or to facilitate the creation of topic-oriented user interfaces that provide the effect of merging unstructured information bases with structured ones.

In short, a *topic map* is a model of knowledge representation based on three key notions: *topics* which represent subjects, *occurrences* of topics which are links to related resources, and *associations* (relations) among topics.

A topic represents, within an application context, any clearly identified and unambiguous subject or concept from the real world: a person, an idea, an object etc.

A topic is an instance of a topic type. Topic types can be structured as hierarchies organized by superclass-subclass relationships. The standard does not provide any predefined semantics to the classes. Finally, topic and topic type form a class-instance relationship.

Topics have three kinds of characteristics: *names* (none, one, or more), *occurrences*, and *roles* in associations. The links between topics and their related information (web page, picture, etc.) are defined by *occurrences*. The linked resources are usually located outside the map. XTM uses a simple link mechanism as defined in XLINK, similar to HTML hyperlinks.

As with topics, occurrences can be typed; occurrence types are often referred to as *occurrence roles*. Occurrence types are also defined as topics. Using XML syntax, the definition of topic is quite simple:

```
<topic id="t-przechlewska-wanda">
  <instanceOf>
    <topicRef xlink:href="#person"/>
  </instanceOf>
  <baseName>
    <baseNameString>Plata-Przechlewska,
      Wanda</baseNameString>
  </baseName>
</topic>
```

Topic associations define relationships between topics. As associations are independent of the resources (i.e., data layer) they represent added-value

information. This independency means that a concrete topic map can describe more than one information pool, and vice versa. Each association can have an *association type* which is also a topic. There are no constraints on how many topics can be related by one association. Topics can play specific roles in associations, described with *association role types* — which are also topics.

The concepts described above are shown in Fig. 1. Topics are represented as small ovals or circles in the upper half of the picture while the large oval at the bottom indicates the data layer. Small objects of different shapes contained in the data layer represent resources of different types. The lines between the data layer and topics represent occurrences, while thick dashed ones between topics depict associations.

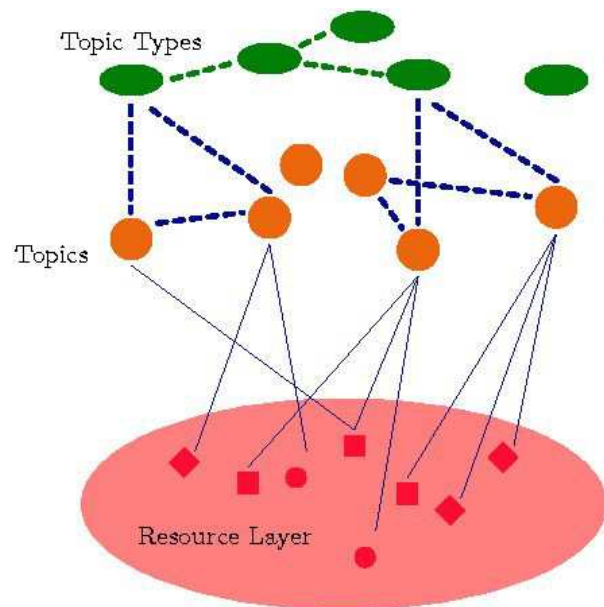


Figure 1: Topic map and resource layer.

Besides the above mentioned three fundamental concepts, the standard provides a notion of *scope*. All characteristics of topics are valid within certain bounds, called a *scope*, and determined in terms of other topics. Typically, scopes are used to model multilingual documents, access rights, different views, and so on.

Scopes can also be used to avoid name conflicts when a single name denotes more than one concept. An example of scope for the topic *latex* might be *computer application* or *rubber industry* depending on the subject of the topic. Only the topic characteristics can be scoped, not the topic itself.

3 Subject Identity and Map Merging

From the above short tour of TM concepts it should be clear that there is an exact one-to-one correspondence between subjects and topics. Thus, the identification of subjects is crucial to individual topic map applications and to interoperability between different topic maps.

The simplest and most popular way of identifying subjects is by identifying them via some system of unique labels (usually URIs). A subject identifier is simply a URI unambiguously identifying the subject. If the subject identifier points to a resource (not required) the resource is called a subject indicator. The subject indicator should contain human-readable documentation describing the non-addressable subject [22].

As there are no restrictions to prevent every map author from defining their own subject identifiers and resource indicators, there is a possibility that semantic or syntactic overlap will occur. To overcome this, *published subject indicators* (PSIs) are proposed [17]. PSIs are stable and reliable indicators published by an institution or organization which desires to promote a specific standard. Anyone can publish PSIs and there is no registration authority. The adoption of PSIs can therefore be an open and spontaneous process [17, 6].²

Subject identity is of primary importance for topic map merging when there is a need to recognize which topics describe the same subject.

Two topics and their characteristics can be *merged* (aggregated) if the topics share the same name in the same scope (*name-based merging*), or if they refer to the same subject indicator (*subject-based merging*). Merging results in a single topic that has the union of all characteristics of merged topics. Merged topics play the roles in all the associations that the individual topics played before [22, 15].

4 Constraining, Querying, and Navigating the Map

The notion of a topic map template is used frequently in literature. As the name suggests, a *topic map template* is a sort of schema imposing constraints on topic map objects with TM syntax. The standard does not provide any means by which the designer of the TM template can put constraints onto the topic map itself. Standardisation of such constraints are currently in progress [14].

² For example, the XTM 1.0 specification contains a set of PSIs for core concepts, such as class, instance, etc., as well as for the identification of countries and languages [19].

Displaying lists of indexes which the user can navigate easily is the standard way of TM visualization. As this approach does not scale well for larger maps, augmenting navigation with some sort of searching facility is recommended. Other visualization techniques such as hyperbolic trees [15], cone trees, and hypergraph views (Fig. 2) can be used for visualization and navigation of topic maps. They display TMs as a graph, with the topics and occurrences as nodes and the associations as arcs. The drawback of such ‘advanced’ techniques is that users are usually unfamiliar with them.

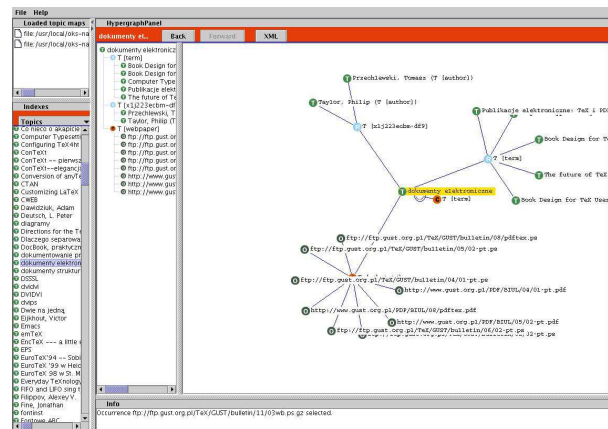


Figure 2: Hypergraph visualization with TMNav.

There are several proposed query languages for topic maps. None of them are part of the standard and there are inconsistencies in different TM engines. Two of the most prominent proposals are:

- TMQL (*Topic Maps Query Language*, [9]), with SQL-like syntax, provides both for querying and modifying topic maps (select, insert, delete, update).
- Tolog, inspired by logic programming language Prolog, supports requirements for TMQL with clearer and simpler syntax.

The introduction to the TM standard presented in this paper does not cover all the details of the technology. Interested readers can find an exhaustive description in [15], which contains detailed introduction with numerous examples, and [16].

5 Topic Maps and RDF

The W3C promotes the Resource Description Framework (RDF) [10] as another framework for expressing metadata. RDF is a W3C standard envisioned to be a foundational layer of the Semantic Web.

The fundamental notion in the RDF data model is a *statement*, which is a triple composed of a

resource, *property*, and *value*. The RDF Schema (RDFS) [4] is a W3C working draft aimed at defining a description language for vocabularies in RDF. More expressive RDFS models have been proposed recently [23].

One key difference between RDF and topic maps is that topic maps are modelled on a concept-centric view of the world. For example, in RDF there are no ‘predefined’ properties, so to assign a name to a resource one has to use another standard (such as Dublin Core), something that is not necessary with topic maps. The notion of scope is also absent from RDF too.

The RDF and Topic Maps standards are similar in many respects [7]. Both offer simple yet powerful means of expressing concepts and relationships.

6 Building Topic Maps for the GUST Bibliographic Database

Similar to writing a good index for a book, creating a good topic map is carried out by combining manual labour with the help of some software applications. It is usually a two-stage task, beginning with the modelling phase of building the ‘upper-part’ of the map, i.e., the hierarchy of topic types and association types (the *schema* of the map) and then populating the map with instances of topic types, their associations and occurrences.

Approaches for developing a topic map out of a pool of information resources include [2]:

- using standard vocabularies and taxonomies (i.e., www.dmoz.org) as the initial source of topic types.
- generating TMs from the structured databases or documents with topic types and association types derived from the scheme of the database/document.
- extraction of topics and topic associations from pools of unstructured or loosely structured documents using NLP (Natural Language Processing) software combined with manual labour.

The first approach is concerned with the modelling phase of topic map creation, while the third one deals with populating the map.

Following the above guidelines, the structure of the BIB_TE_X records was an obvious candidate to start with in modelling our map of GUST articles. It provides a basic initial set of topics including: *author*, *paper*, *keyword*, and the following association types: *author-paper*, *paper-keyword* and *author-keyword*. Abstracts (if present in BIB_TE_X databases) can be considered as occurrences of the topic *paper*.

The publication date and language can be used as scopes for easy navigation using them as constraints.

Other TAOs (topics, associations, and occurrences [16]) to consider are: author home pages (occurrence type), applications described in the paper (association type), papers referenced (association type). This information is absent from BIB_TE_X files but, at least theoretically, can be automatically extracted from the source files of papers.

We started by analyzing the data at our disposal, i.e., T_EX and BIB_TE_X source files. Unfortunately, in the case of the GUST bulletin the BIB_TE_X database was not maintained. This apparent oversight was rectified with simple Perl scripts and a few days of manual labour. The bibliographic database was created and saved in a XML-compatible file.³

T_EX documents are often visually tagged and lack information oriented markup. The only elements marked up consistently and unambiguously in the case of the GUST bulletin are the paper titles and authors’ names. Authors’ home pages were rarely present, while email addresses were available but not particularly useful for our purposes. Neither abstracts nor keyword lists had been required and as a consequence were absent from the majority of the papers. Similarly, any consistent scheme of marking bibliographies (or attaching *.bib* files) was lacking, so there was no easy way to define the *related to* association between papers.

The benefit derived from keywords is much greater if they are applied consistently according to some fixed classification; otherwise, the set of keywords usually consists of many random terms which are nearly useless. Since we didn’t want to define yet another ‘standard’ in this area, we would have liked to adopt an existing one. The following sources were considered: the T_EX entry at [dmoz.org](http://www.dmoz.org), Graham Williams’ catalogue.⁴, collections of BIB_TE_X files and *.tpm* files [20]

The accuracy of the T_EX taxonomy subtree at [dmoz.org](http://www.dmoz.org) was somewhat questionable, and we quickly rejected the idea of using it. Williams’ catalogue of T_EX resources does not include any information except the location of the resource in the structure of CTAN. As for BIB_TE_X files, it appeared only MAPS and *TUGboat* were complete and up-to-date⁵ but only the latter contains keywords. Unfortunately, they don’t comply with any consistent

³ We reused the XML schema developed for the MAPS bulletin (<http://www.ntg.org.ln/maps/>).

⁴ <http://www.ctan.org/tex-archive/help/Catalogue>

⁵ Cahiers GUTenberg was not found, but the impressive portal of Association GUTenberg indicates appropriate meta-data are maintained, but not published.

scheme. Due to the lack of any existing standard, the keywords were added manually on a common-sense basis, with the intention of being ‘in sync’ with the most frequent terms used in MAPS.⁶

Finally the following TAOs were defined (the language of the publication was considered to be the only scope):

- topic types: *author*, *paper*, and *keyword*;
- association types: *author-paper*, *paper-keyword*, and *author-keyword*;
- occurrence types: *papers* and *abstracts*.

The schema of the map was prepared manually and then the rest of the map was generated from the content of intermediate XML file with an XSLT stylesheet [8, 13]. The resulting map consists of 454 topics, 1029 associations, and 999 occurrences. A fragment of the map rendered in a web browser with Ontopia Omnigator (a no-cost but closed-source application, <http://www.ontopia.net/download/>) is shown in Fig. 3.

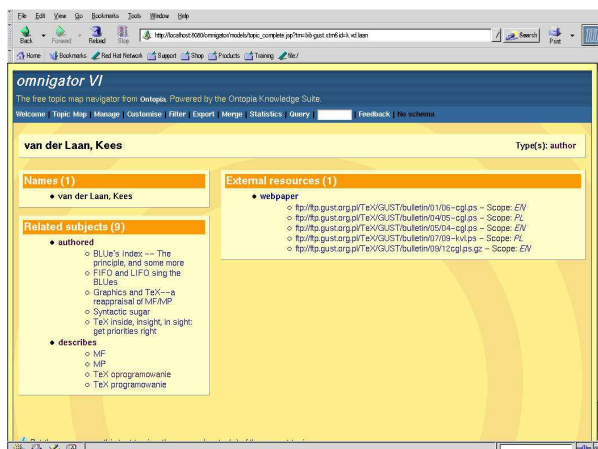


Figure 3: A fragment of GUST topic map rendered with Omnigator.

Omnigator shows the map as a list of links arranged in panels. Initially only a list of subjects (index of topic types) is displayed. When a link for a topic is clicked on, the topic is displayed with all the information about its characteristics (names, associations, occurrences). In Fig. 3, an example page for author *Kees van der Laan* is shown. The right panel contains all the relevant resources while the

⁶ There are 814 bibliographic entries in MAPS base and 895 different keywords. The most popular keywords in MAPS BibT_EX file are: L^AT_EX – 51, NTG – 42, plain T_EX – 37, PostScript – 28, ConT_EXt, T_EX-NL, METAFONT, SGML, and so forth. There are small number of inconsistent cases (special commands vs. specials, or configuration vs. configuring) and fine-grained keywords (Portland, Poland, Bachotek, USSR!).

lower left has all the related topics, i.e., papers written by Kees and other subjects covered. The user can easily browse both papers authored by him and switch to pages on some other interesting subject. The panel with resources contains information on the resource type allowing fast access to the required data.

Similar functionality can be obtained with the freely available TM4Nav or even by using a simple XSLT stylesheet [13].

7 Integrating Other T_EX Resources

So far there is nothing in TMs which cannot be obtained using other technologies. The same or better functionality can be achieved with any database management system (DMS). But integrating T_EX resources on a global scale needs flexibility, which traditional RDBMS-based DMS applications lack. For example, topic maps can be extended easily through merging separate maps into one, while DMS-based extensions usually require some prior agreement between the parties (e.g., LUGs), schema redefinitions, and more.

To verify this flexibility in practice, we extended the GUST map with the MAPS and TUB BibT_EX databases. For easy interoperability in a multi-language environment, the upper half of the map was transferred to a separate file. With the use of scope, the design of multi-language topic types was easy, for example:

```
<topic id="english">
  <subjectIdentity>
    <subjectIndicatorRef
      xlink:href="http://www.topicmaps.org/\
        xtm/1.0/language.xtm#en"/>
    </subjectIdentity>
    <baseName>
      <baseNameString>EN</baseNameString>
    </baseName>
  </topic>
  ...
  <topic id="author">
    <baseName><scope>
      <topicRef xlink:href="#english"/></scope>
      <baseNameString>author</baseNameString>
    </baseName>
    <baseName><scope>
      <topicRef xlink:href="#polish"/> </scope>
      <baseNameString>autor</baseNameString>
    </baseName>
  </topic>
```

Other topic types were designed similarly. Scopes for other languages can easily be added.

The ‘lower part’ of the map was generated from (cleaned) BibT_EX records with `bibtex2xml.py`

(<http://bibtexml.sf.net>) and than transformed to MAPS XML with an XSLT stylesheet. Keywords were added to TUB entries using a very crude procedure.⁷

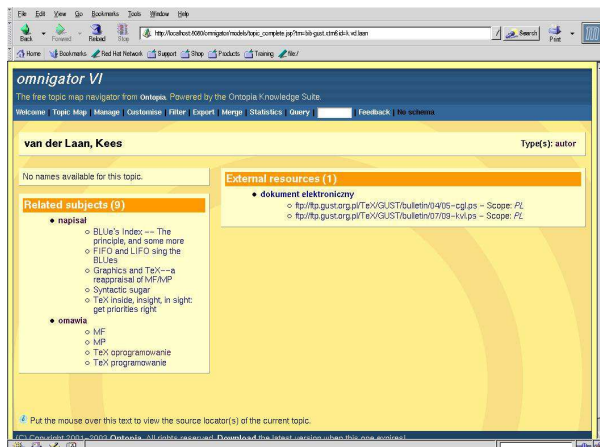


Figure 4: Topic map fragment from Fig. 3 scoped to Polish language.

Simple name-based merging of all three maps results in over 25,000 TAOs (≈ 1000 authors, more than 2000 papers). Some of the subjects were represented with multiple topics. As an example the Grand Wizard was represented as the following four distinct topics: ‘Knuth, Don’, ‘Knuth, Donald’, ‘Knuth, Donald E.’, ‘Knuth., Donald E.’⁸

As identity-based merging is regarded as more robust, some identifiers have to be devised first. Establishing a PSI for every \TeX author seemed overly ambitious. Instead, a dummy subject identifier was chosen, such as: <http://tug.org/authors#initials-surname>. This can still produce multiple topics for the same subject, but now we can eliminate unwanted duplicates by defining an additional map consisting solely of topics like the following [18]:

```
<topic id="de-knuth">
  <subjectIdentity>
    <subjectIndicatorRef
      xlink="http://tug.org/authors#d-knuth"/>
    <subjectIndicatorRef
      xlink="http://tug.org/authors#d-e-knuth"/>
  </subjectIdentity>
</topic>
```

Merging this map with the ‘base’ map(s) will result in a map free of unwanted duplicated topics with all variant names preserved.

⁷ Acronyms, such as \LaTeX , METAFONT, or XML, present in the title were used as keywords.

⁸ First name variants, abbreviations and middle names cause problems in many more cases.

For further extensions, we plan to incorporate CTAN resources. For that purpose, Williams’ catalogue and/or the TPM files from \TeX Live project can be used. As the catalogue contains author names, it would be for example possible to enrich the map with the *author-application* association. Further enrichment will result if we can link applications with documents describing them. However, some robust classification scheme of \TeX resources should be devised first.

8 Topic Map Tools

As with any other XML-based technology, topic maps can be developed with any text editor and processed with many XML tools. However, for larger-scale projects specialized software is needed. There are a few tools supporting topic map technology, developed both commercially and as Open Source projects. We have considered both Ontopia Omnigator (mentioned in the previous section) and TM4J (free software).

TM4J (<http://tm4j.org>) is a suite of Java packages which provide interfaces and default implementations for the import, manipulation and export of XML Topic Maps. Features of the TM4J engine include an object model which supports XTM specification with the ability to store topic map in an object-oriented or relational database, and an implementation of the tolog query language.

Based on TM4J a few projects are in progress: TMNav for intuitive navigation and editing of topic maps, and TMBrowse for publishing maps as set of HTML pages (similarly to Omnigator).

These projects are in early stages and our experience with TMBrowse indicates that current version frequently crashes with bigger maps and is significantly slower than Omnigator. There were problems with tolog queries as well.

As all these projects are actively maintained progress may be expected in the near future.

9 Summary

Topic maps are an interesting new technology which can be used to describe the relation between \TeX resources. The main problem is topic map visualization. Available tools are in many cases unstable and non-scalable, but we can expect improvement.

The system presented here can certainly be improved. It is planned to extend it with the content of Williams’ catalogue. The maps developed in the project are available from <http://gnu.univ.gda.pl/~tomasz/tm/>. At the same address, the interested reader can find links to many resources on topic maps.

References

- [1] ISO/IEC. Topic Maps, ISO/IEC 13250, 2002. <http://www.y12.doe.gov/sgml/>.
- [2] Kal Ahmed, Danny Ayers, Mark Birbeck, and Jay Cousins. *Professional XML Meta Data*. Wrox Press, 2001.
- [3] Tim Berners-Lee, James Hendler, and Ora Lassila. The semantic web. *Scientific American*, 284(5):35–43, 2001. <http://www.sciam.com/2001/0501issue/0501berners-lee.html>.
- [4] Dan Brickley and R. V. Guha. RDF vocabulary description language 1.0: RDF schema, 2002. <http://www.w3.org/TR/rdf-schema/>.
- [5] Anna Carlstedt and Mats Nordborg. An evaluation of topic maps. Master's thesis, Göteborg University, 2002. <http://www.cling.gu.se/~c18matsn/exjobb.html>.
- [6] Paolo Ciancarini, Marco Pirruccio, and Fabio Vitali. Metadata on the Web. On the integration of RDF and topic maps. In *Extreme Markup Languages*, 2003. <http://www.mulberrytech.com/Extreme/Proceedings/xslfo-pdf/2003/Presutti01/EML2003Presutti01.pdf>.
- [7] Lars M. Garshol. Living with Topic Maps and RDF. <http://www.ontopia.net/topicmaps/materials/tmrdf.html>, 2002.
- [8] Michael Kay. *XSLT Programmer's Reference 2nd Edition*. Wrox Press, 2001.
- [9] Rafał Ksieżyk. Trying not to get lost with a topic map. In *XML Europe Conference*, 1999. <http://www.infoloom.com/gcaconfs/WEB/granada99/ksi.HTM>.
- [10] Ora Lassila and Ralph R. Swick. Resource description framework (RDF). Model and syntax specification, 1999. <http://www.w3.org/TR/REC-rdf-syntax/>.
- [11] Xia Lin and Jian Qin. Building a topic map repository, 2000. <http://www.knowledgetechnologies.net/proceedings/presentations/lin/xialin.pdf>.
- [12] Ashish Mahabal, S. George Djorgovski, Robert Brunner, and Roy Williams. Topic maps as a virtual observatory tool, 2002. <http://arxiv.org/abs/astro-ph/0110184>.
- [13] Sal Mangano. *XSLT Cookbook*. O'Reilly, 2002.
- [14] Mary Nishikawa and Graham Moore. Topic map constraint language requirements, 2002. <http://www.isotopicmaps.org/tmcl/>.
- [15] Jack Park and Sam Hunting, editors. *XML Topic Maps. Creating and using Topic Maps for the Web*. Addison-Wesley, 2002.
- [16] Steve Pepper. The TAO of topic maps, 2000. <http://www.ontopia.net/topicmaps/materials/tao.html>.
- [17] Steve Pepper. Published subject: Introduction and basic requirements, 2003. <http://www.oasis-open.org/committees/>.
- [18] Steve Pepper and Marius L. Garshol. Lessons on applying topic maps. <http://www.ontopia.net/topicmaps/materials/xmlconf.html>, 2002.
- [19] Steve Pepper and Graham Moore. XML topic maps (XTM) 1.0, 2000. <http://www.topicmaps.org/xtm/1.0/>.
- [20] Fabrice Popineau. Directions for the TeXlive systems. In *EuroT_EX 2001, The Good, the Bad and the Ugly, Kerkrade*, pages 152–161. NTG, 2001. http://www.ntg.nl/maps/pdf/26_20.pdf.
- [21] Tomasz Przechlewski. Wykorzystanie map pojęć w zarządzaniu repozytoriami dokumentów elektronicznych. In *Materiały Konferencji: MSK 2003*, 2003.
- [22] Hans Holger Rath. Semantic resource exploitation with topic maps. In *Proceedings of the GLDV-Spring Meeting 2001*, 2001. <http://www.uni-giessen.de/fb09/ascl/gldv2001/>.
- [23] Michael K. Smith, Chris Welty, and Deborah L. McGuinness. OWL web ontology language guide, 2003. <http://www.w3.org/TR/owl-guide/>.

Interactive Editing of MathML Markup Using \TeX Syntax*

Luca Padovani

Department of Computer Science, University of Bologna

Mura Anteo Zamboni, 7

40127 Bologna

Italy

lpadovan@cs.unibo.it

<http://www.cs.unibo.it/~lpadovan/>

Abstract

We describe the architecture of a syntax-directed editor for authoring structured mathematical documents which can be used for the generation of MathML markup [4]. The author interacts with the editor by typing \TeX markup as he would do in a normal text editor, with the difference that the typed markup is parsed and displayed on-the-fly. We discuss issues regarding both the parsing and presentation phases and we propose implementations for them. In contrast with existing similar tools, the architecture we propose offers better compatibility with \TeX syntax, a pervasive use of standard technologies and a clearer separation of content and presentation aspects of the information.

1 Introduction

MathML [4] is an XML [2] application for the representation of mathematical expressions. Like most XML applications, MathML is unsuitable to be written directly because of its verbosity except in the simplest cases. Hence the editing of MathML documents needs the assistance of dedicated tools. As of today, such tools can be classified into two main categories:

1. WYSIWYG (What You See Is What You Get) editors that allow the author to see the formatted document on the screen while it is being composed. The editor usually provides some “export mechanism” that creates XML with embedded MathML from the internal representation of the document;
2. conversion tools that generate MathML markup from different sources, typically other markup languages for scientific documents, such as \TeX [5].

Tools in the first category are appealing, but they suffer from at least two limitations: a) editing is typically *presentation oriented* — the author is primarily concerned about the “look” of the document and tends to forget about its content. b) They may slow down the editing process because they often involve the use of menus, palettes of symbols,

and, in general, the pointing device for completing most operations.

In this paper we describe the architecture of a tool that tries to synthesize the “best of both worlds”. The basic idea is to create a WYSIWYG editor in which editing is achieved by typing concrete markup as the author would do in an actual plain text editor. The markup is then tokenized and parsed on-the-fly, a corresponding presentation is created by means of suitable transformations, and finally displayed. The editor is meant not only as an authoring tool, but more generally as an interface for math applications.

Although in the paper we assume that the concrete markup typed by the user is \TeX (more precisely the subset of \TeX concerned about mathematics) and that presentation markup is MathML, the system we are presenting is by no means tied to these languages and can be targeted to other contexts as well. One question that could arise is: “why \TeX syntax?” We can see at least three motivations: first of all because of \TeX popularity in many communities. Second, because macros, which are a fundamental concept in \TeX , are also the key to editing at a more content-oriented level, which is a primary requirement for many applications handling mathematics. Finally, because, as we will see, \TeX markup has good *locality* properties which make it suitable in the interactive environment of our concern.

* This work has been supported by the European Project IST-2001-33562 MoWGLI.

The body of the paper is structured into four main sections: in Section 2 we overview the architecture of the tool while in Sections 3, 4, 5 we describe in more detail the main phases of the editing process (lexing, parsing, and transformation). Familiarity with T_EX syntax and XML-related technologies is assumed.

2 Architecture

Several tools for the conversion of T_EX markup suffer from two major drawbacks that we are not willing to tolerate in our design: (1) they rely on the T_EX system itself for parsing the markup. While guaranteeing perfect compatibility with T_EX, this implies the installation of the whole system. Moreover, the original T_EX parser does not meet the incremental requirements that we need; (2) the lack of flexibility in the generation of the target document representation, which is either fixed by the conversion tool or it is only slightly customizable by the user.

To cope with problem (1) we need to write our own parser for T_EX markup. This is well known to be a non-trivial task, because of some fancy aspects regarding the very nature of T_EX syntax and the lack of a proper “T_EX grammar”. We will commit ourselves with a subset of T_EX syntax which appears to be just what an average author needs when writing a document. As we will see, the loss in the range of syntactic expression is compensated by a cleaner and more general transformation phase. As for the lack of a T_EX grammar, we perceive this as a feature rather than a weakness: after all T_EX is built around the fact that authors are free to define their own macros. Macros are the fundamental entities giving structure to the document.

Let us now turn our attention to problem (2): recall that the general form of a T_EX macro definition (see *The T_EXbook*, [5]) is

```
\def⟨control sequence⟩⟨parameter text⟩
  {⟨replacement text⟩}
```

where the ⟨parameter text⟩ gives the syntax for invoking the macro and its parameters whereas the ⟨replacement text⟩ defines somehow the “semantics” of the macro (typically a presentational semantics). Thus the ultimate semantic load of a macro is invariably associated with the configuration of the macro at the point of definition.

We solve problem (2) by splitting up macro definitions so that structure and semantics can be treated independently. A well-formed T_EX document can be represented as a tree whose leaves are either literals (strings of characters) or macros with no parameters, and each internal node represents a

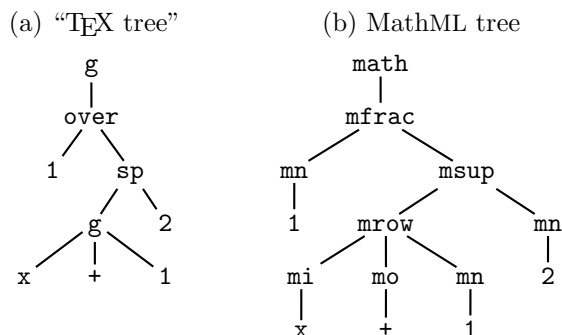


Figure 1: Tree representation for $\{1\over{x+1}\}^2$ and corresponding MathML markup.

macro and the node’s children are the macro’s parameters. Entities like delimiters, square brackets surrounding optional parameters or literals occurring in the ⟨parameter text⟩ of macro definitions are purely syntactic and need not be represented in the tree if our main concern is capturing the structure of the document. Fig. 1(a) shows the tree structure of a simple mathematical formula.

Once the document is represented as a tree, the process of macro expansion — that is, *interpretation* — can be defined as a recursive transformation on the nodes of the tree. As we will represent trees using XML, transformations can be very naturally implemented by means of XSLT stylesheets [3]. Fig. 1(b) shows the MathML tree corresponding to the T_EX tree on the left hand side. The two trees are basically isomorphic except for the name of the nodes and the presence of explicit token nodes for literals in the MathML tree. This is to say that the MathML tree can be generated from the T_EX tree by simple transformations. However, once the interpretation phase is independent of parsing (something which does not happen in T_EX) it is natural to define much more general transformations that are not just node-by-node rewritings.

The following are the main components of an interactive, syntax-based editor for structured documents:

INPUT BUFFER: the sequence of concrete characters typed by the author;

LEXICAL ANALYZER: responsible for the tokenization of the characters in the input buffer;

DICTIONARY: a map from ⟨control sequence⟩ to ⟨parameter text⟩ which is used to know the syntax of macros;

PARSER: for the creation of the internal tree structure representing the document;

TRANSFORMATION ENGINE: for mapping the internal tree into the desired format.

No doubt these entities are common to all tools converting \TeX markup into a different format, but the degree of mutual interdependence and the way they are implemented may differ considerably, especially when interactivity is a main concern. The added value of our approach is that it allows the author to independently customize both the dictionary and the transformation engine, and the advanced user of the editor the possibility of adapting the lexical analyzer to languages other than \TeX (we will spend a few more words on this topic in the conclusions).

Notation We will use the following conventions regarding lists. Lists are uniformly typed, that is elements of a list are all of the same type. We use α^* to denote the type of a list whose elements have type α . \square is the empty list; $n :: x$ is the list with head element n and tail x ; $x@y$ is the concatenation of two lists x and y ; $[n_1; \dots; n_k]$ is a short form for $n_1 :: \dots :: n_k :: \square$.

3 Lexical Analysis

The purpose of this phase is to tokenize the input buffer. As we are talking about an interactive tool, the presence of an input buffer may look surprising. Implementations for the input buffer range from *virtual* buffers (there is no buffer at all, characters are collected by the lexical analyzer which outputs tokens as they are completed) to *flat* buffers (just a string of characters as in a text editor) to *structured* buffers. For efficiency, we do not investigate in detail all the possibilities in this paper, but early experiments have shown that working with virtual buffers can be extremely difficult. As long as insert operations are performed at the right end of the buffer the restructuring operations on the parsing tree are fairly easy, but when it comes to deletion or to modifications in arbitrary positions, the complexity of restructuring operations rises rapidly to an unmanageable level. Hence, from now on we will assume that a flat input buffer is available. Whether the buffer should be visible or not is a subjective matter, and may also depend on the kind of visual feedback given by the editor on incomplete and/or incorrect typed markup.

The outcome of the lexer is a stream (list) of tokens. Each token may have one of three forms: a *literal*, that is a single character to be treated “as is”, a *space*, that is a sequence of one or more space-like characters, or a *control sequence*, that is the name of a macro.

Since the token stream is the only interface between the lexer and the parser, the lexer has the freedom to perform arbitrary mappings from the characters in the input buffer to tokens in the stream. In particular, some \TeX commands like $\backslash\alpha$ or $\backslash\rightarrow$ are just placeholders for Unicode characters. There is no point in communicating these entities as control sequences as the internal tree representation (XML) is able to accommodate Unicode characters naturally; also, treating them as literals simplifies the subsequent transformation phase.

On the other hand, there are characters, such as curly braces $\{$ and $\}$ or scripting operators $_$ and \hat , that have a special meaning. Logically these are just short names for macros that obey their own rules regarding parameters. What we propose is a general classification of parameter types which, in addition to parameters in normal \TeX definitions, allows us

- to deal with optional parameters as \LaTeX [6] does;
- to treat $\{$ as just an abbreviation for $\backslash\text{bgroup}$ and make $\backslash\text{bgroup}$ a macro with one parameter delimited by $\backslash\text{egroup}$, which we treat as the expansion for $\}$. In order for this “trick” to work we have to design the parser carefully, as we will see in Sect. 4;
- to treat scripting operators $_$ and \hat as the two macros $\backslash\text{sb}$ and $\backslash\text{sp}$ both accepting a so-called pre-parameter (a parameter that occurs *before* the macro in the input buffer) and a so-called post-parameter (a parameter that occurs *after* the macro in the input buffer);
- to deal with macros that have “open” parameters. For instance $\backslash\text{rm}$ affects the markup following it until the first delimiter coming from an outermost macro is met. We treat $\backslash\text{rm}$ as a macro with an open post-parameter that extends as far as possible to the right. Similarly, $\backslash\text{over}$ can be seen as a macro with open pre- and post-parameters.

In order to describe parameter types we need to define the concept of *term*. A *term* is either a literal or a macro along with all its parameters (equivalently, a term is a subtree in the parsing tree). A *simple* parameter consists of one sole term. A *compound* parameter consists of one or more terms extending as far as possible to the left or to the right of the macro depending on whether the parameter is “pre-” or “post-”. A *delimited* parameter consists of one or more terms extending as far as possible to the right up to but not including a given token

Table 1: Examples of T_EX and L^AT_EX macros along with their signature.

Macro	Parameters	
	pre	post
overline		[simple]
sqrt		[simple] (T _E X)
		[optional; simple] (L ^A T _E X)
root		[delimited(control(of)); simple]
over, choose	[compound]	[compound]
frac		[simple; simple]
rm, bf, tt, it		[compound]
left		[simple; delimited(control(right)); simple]
sb, sp	[simple]	[simple]
bggroup		[delimited(control(egroup))]
begin		[simple; optional; delimited(control(end)); simple]
proclaim		[token(space); delimited(literal(.)); token(space); delimited(control(par))]

t. An *optional* parameter is either empty or it consists of one or more terms enclosed within a pair of square brackets [and]. The absence of the opening bracket means that the optional parameter is not given. A *token* parameter is a given token *t* representing pure syntactic sugar. It does not properly qualify as a parameter and does not appear in the parsing tree.

Formally tokens and parameter types are defined as follows:

$$\begin{aligned}
 \textit{token} &::= \textit{literal}(v) \mid \textit{space} \mid \textit{control}_{(p_1, p_2)}(v) \\
 \textit{type} &::= \textit{simple} \mid \textit{compound} \mid \textit{delimited}(t) \\
 &\quad \mid \textit{optional} \mid \textit{token}(t)
 \end{aligned}$$

where $t \in \textit{token}$, $v \in \textit{string}$ is an arbitrary string of Unicode characters, $p_1 \in \{\textit{simple}, \textit{compound}\}^*$ and $p_2 \in \textit{type}^*$ are lists of parameter types for the pre- and post-parameters respectively. Note that pre-parameters can be of type *simple* or *compound* only.

The dictionary is a total map

$$\textit{dictionary} : \textit{string} \mapsto \textit{token}$$

such that for each unknown control sequence v we have $\textit{dictionary}(v) = \textit{control}_{\langle [], [] \rangle}(v)$. Table 1 shows part of a possible dictionary for some T_EX and L^AT_EX commands (mostly for mathematics). Note how it is possible to encode the signature for the `\begin` control sequence, although it is not possible to enforce the constraint that the first and the last parameters must have equal value in order for the construct to be balanced.

4 Parsing

We now come to the problem of building the T_EX parsing tree starting the stream of tokens produced by the lexical analyzer. As we have already pointed

out there is no fixed grammar that we can use to generate the parser automatically: authors are free to introduce new macros and hence new ways of structuring the parse tree. Thus we will build the parser “by hand”. More reasons for writing an ad-hoc parser, namely error recovery and incrementality, will be discussed later in this section.

The following grammar captures formally the structure of a T_EX parsing tree, which is the outcome of the parser:

$$\begin{aligned}
 \textit{node} &::= \textit{empty} \\
 &\quad \mid \textit{literal}(v) \quad v \in \textit{string} \\
 &\quad \mid \textit{macro}(v, x) \quad v \in \textit{string}, x \in \textit{param}^* \\
 \textit{param} &::= \{a\} \quad a \in \textit{node}^*
 \end{aligned}$$

Note that a parameter is made of a list of nodes and that literals are strings instead of single characters. The `empty` node is used to denote a missing term when one was expected; its role will be clarified later in this section.

The appendix contains the Document Type Definition for the XML representation of T_EX parsing trees. It is simpler than the T_EXML DTD [7] and we are providing it as mere reference.

4.1 Parsing Functions

Table 2 gives the operational semantics of the parser. In this table only, for each $a \in \textit{node}^*$ we define $a! = [\textit{empty}]$ if $a = []$ and $a! = a$ otherwise. There are four parsing functions: \mathcal{T} for terms, \mathcal{A} for pre-parameters, \mathcal{B} for post-parameters, and \mathcal{C} for delimited sequences of terms. Each parsing function is defined by induction on the structure of its arguments. Axioms (rules with no horizontal line) denote base cases, while inference rules define the value

Table 2: Parsing functions for the simplified \TeX markup.

$\forall d \in \text{token}^*$	$\xrightarrow{\mathcal{T}(d)}$	$:$	$\text{node}^* \times \text{token}^* \rightarrow \text{node}^* \times \text{token}^*$
	$\xrightarrow{\mathcal{A}}$	$:$	$\text{node}^* \times \text{type}^* \rightarrow \text{node}^* \times \text{param}^*$
$\forall d \in \text{token}^*$	$\xrightarrow{\mathcal{B}(d)}$	$:$	$\text{type}^* \times \text{token}^* \rightarrow \text{param}^* \times \text{token}^*$
$\forall d \in \text{token}^*, \forall b \in \text{bool}$	$\xrightarrow{\mathcal{C}(d,b)}$	$:$	$\text{node}^* \times \text{token}^* \rightarrow \text{node}^* \times \text{token}^*$

$$\begin{array}{l}
\text{(T.1)} \ a, [] \xrightarrow{\mathcal{T}(d)} a, [] \quad \text{(T.2)} \ a, t :: l \xrightarrow{\mathcal{T}(d)} a, t :: l \quad (t \text{ occurs in } d) \\
\text{(T.3)} \ a, \text{literal}(v) :: l \xrightarrow{\mathcal{T}(d)} a@[\text{literal}(v)], l \quad \text{(T.4)} \ \frac{a, l \xrightarrow{\mathcal{T}(d)} a', l'}{a, \text{space} :: l \xrightarrow{\mathcal{T}(d)} a', l'} \\
\text{(T.5)} \ \frac{a, p_1 \xrightarrow{\mathcal{A}} a', x \quad p_2, l \xrightarrow{\mathcal{B}(d)} y, l'}{a, \text{control}_{\langle p_1, p_2 \rangle}(v) :: l \xrightarrow{\mathcal{T}(d)} a'@[\text{macro}(v, x@y)], l'} \\
\text{(A.1)} \ a, [] \xrightarrow{\mathcal{A}} a, [] \quad \text{(A.2)} \ \frac{[], p \xrightarrow{\mathcal{A}} a, x}{[], s :: p \xrightarrow{\mathcal{A}} a, x@[\{\{\text{empty}\}\}]} \\
\text{(A.3)} \ \frac{a, p \xrightarrow{\mathcal{A}} a', x}{a@[n], \text{simple} :: p \xrightarrow{\mathcal{A}} a', x@[\{\{n\}\}]} \quad \text{(A.4)} \ \frac{[], p \xrightarrow{\mathcal{A}} a', x}{a, \text{compound} :: p \xrightarrow{\mathcal{A}} a', x@[\{\{a\}\}]} \\
\text{(B.1)} \ [], l \xrightarrow{\mathcal{B}(d)} [], l \quad \text{(B.2)} \ \frac{p, t' :: l \xrightarrow{\mathcal{B}(d)} x, a}{\text{token}(t) :: p, t' :: l \xrightarrow{\mathcal{B}(d)} x, a} \dagger \quad (t \neq t') \\
\text{(B.3)} \ \frac{p, [] \xrightarrow{\mathcal{B}(d)} x, l}{\text{token}(t) :: p, [] \xrightarrow{\mathcal{B}(d)} x, l} \dagger \quad \text{(B.4)} \ \frac{p, l \xrightarrow{\mathcal{B}(d)} x, l'}{\text{token}(t) :: p, t :: l \xrightarrow{\mathcal{B}(d)} x, l'} \\
\text{(B.5)} \ \frac{[], l \xrightarrow{\mathcal{T}(d)} a, l' \quad p, l' \xrightarrow{\mathcal{B}(d)} x, l''}{\text{simple} :: p, l \xrightarrow{\mathcal{B}(d)} \{a!\} :: x, l''} \quad \text{(B.6)} \ \frac{[], l \xrightarrow{\mathcal{C}(d, \text{false})} a, l' \quad p, l' \xrightarrow{\mathcal{B}(d)} x, l''}{\text{compound} :: p, l \xrightarrow{\mathcal{B}(d)} \{a!\} :: x, l''} \\
\text{(B.7)} \ \frac{p, [] \xrightarrow{\mathcal{B}(d)} x, l}{\text{optional} :: p, [] \xrightarrow{\mathcal{B}(d)} \{\}\} :: x, l \quad \text{(B.8)} \ \frac{[], l \xrightarrow{\mathcal{C}(\text{literal}(\cdot)) :: d, \text{true}} a, l' \quad p, l' \xrightarrow{\mathcal{B}(d)} x, l''}{\text{optional} :: p, \text{literal}(\cdot) :: l \xrightarrow{\mathcal{B}(d)} \{a\} :: x, l''} \\
\text{(B.9)} \ \frac{p, t :: l \xrightarrow{\mathcal{B}(d)} x, l'}{\text{optional} :: p, t :: l \xrightarrow{\mathcal{B}(d)} \{\}\} :: x, l' \quad (t \neq \text{literal}(\cdot)) \\
\text{(B.10)} \ \frac{[], l \xrightarrow{\mathcal{C}(t :: d, \text{true})} a, l' \quad p, l' \xrightarrow{\mathcal{B}(d)} x, l''}{\text{delimited}(t) :: p, l \xrightarrow{\mathcal{B}(d)} \{a!\} :: x, l''} \\
\text{(C.1)} \ a, [] \xrightarrow{\mathcal{C}(d,b)} a, [] \\
\text{(C.2)} \ a, t :: l \xrightarrow{\mathcal{C}(t :: d, \text{true})} a, l \quad \text{(C.3)} \ a, t :: l \xrightarrow{\mathcal{C}(d,b)} a, t :: l \quad (t \text{ occurs in } d) \\
\text{(C.4)} \ \frac{a, t :: l \xrightarrow{\mathcal{T}(d)} a', l' \quad a', l' \xrightarrow{\mathcal{C}(d,b)} a'', l''}{a, t :: l \xrightarrow{\mathcal{C}(d,b)} a'', l''} \quad (t \notin d)
\end{array}$$

of a parsing function (the conclusion, below the line) in terms of the value of one or more recursive calls to other functions (the premises, above the line). Right arrows denote the action of parsing. Arrows are decorated with a label that identifies the parser along with its parameters, if any. The \mathcal{T} , \mathcal{B} , and \mathcal{C} parsers have a parameter representing the list of delimiters in the order they are expected, with the head of the list being the first expected delimiter. The \mathcal{C} parser also has a boolean parameter indicating whether the parser should or should not “eat” the delimiter when it is eventually met.

The root parsing function is \mathcal{T} . Given a delimiter $t \in \text{token}$ and a token stream $l \in \text{token}^*$ we have

$$\square, l \xrightarrow{\mathcal{T}([t])} [n], l'$$

where $n \in \text{node}$ is the parsed term and $l' \in \text{token}^*$ is the part of the token stream that has not been consumed. Spaces are ignored when parsing terms and pre-parameters (rule $T.4$), but not when parsing post-parameters (rule $B.4$). The \mathcal{A} function differs from the other parsing functions because by the time a macro with pre-parameters is encountered, pre-parameters have already been parsed. The lists $a \in \text{node}^*$ in the \mathcal{T} , \mathcal{A} , and \mathcal{C} parsers represent the terms accumulated before the term being parsed. Note that pre-parameters are inserted at the end of the parameter list (rules $A.2$ to $A.4$) and that post-parameters are inserted at the beginning of the parameter list (rules $B.5$ to $B.10$). This way parameter nodes appear in the parse tree in the same order as in the original token stream (rule $T.5$).

4.1.1 Example

Given that the input buffer contains the TeX source shown in Fig. 1, the lexical analyzer would produce the following stream of tokens:

$$l_0 \stackrel{\text{def}}{=} [\text{control}_{\langle \square, [\text{delimited}(\text{control}(\text{egroup})) \rangle]}(\text{bgroup}); \\ \text{literal}(1); \text{control}_{\langle [\text{compound}], [\text{compound}] \rangle}(\text{over}); \\ \text{control}_{\langle \square, [\text{delimited}(\text{control}(\text{egroup})) \rangle]}(\text{bgroup}); \\ \text{literal}(x); \text{literal}(+); \text{literal}(1); \\ \text{control}(\text{egroup}); \text{control}_{\langle [\text{simple}], [\text{simple}] \rangle}(\text{sp}); \\ \text{literal}(2); \text{control}(\text{egroup})]$$

By the application of the parsing rules given in Table 2 it can be shown that

$$\square, l_0 @ [\text{control}(\text{eoi})] \xrightarrow{\mathcal{T}([\text{control}(\text{eoi})])} [n], [\text{control}(\text{eoi})]$$

where $n \in \text{node}$ is the same tree shown in Fig. 1 except that the g nodes are labeled with bgroup .

4.2 Error Recovery

Parsing functions are all total functions, they always produce a result, even when the input token

stream is malformed. Unlike parsers of batch TeX converters or the TeX parser itself, there will often be moments during the editing process when the input buffer contains incorrect or incomplete markup, for example because not all the required parameters of a macro have been entered yet. The parser must recover from such situations in a tolerant and hopefully sensible way. We distinguish three kinds of situations: *missing parameters*, *pattern mismatch*, and *ambiguity*, which we examine in the rest of this section.

4.2.1 Missing Parameters

Consider an input token stream representing the sole $\overline{\hspace{1cm}}$ macro with no arguments provided:

$$l_1 \stackrel{\text{def}}{=} [\text{control}_{\langle [\text{compound}], [\text{compound}] \rangle}(\text{over}); \\ \text{control}(\text{eoi})]$$

It is easy to check that

$$\square, l_1 \xrightarrow{\mathcal{T}([\text{control}(\text{eoi})])} [\text{macro}(\text{over}, [\text{empty}; \text{empty}]), \\ [\text{control}(\text{eoi})]]$$

More generally the parser inserts `empty` nodes in the parsing tree wherever an expected parameter is not found in the token stream. This behavior can be seen in rule $A.2$ and also in rules $B.5$, $B.6$, and $B.10$ where the `!` operator is used. For optional parameters an empty node list is admitted (rules $B.7$ and $B.8$).

The presence of `empty` nodes guarantees that the generated tree is structurally well-formed, which is crucial for the subsequent transformation phase. It also allows the application to give the user feedback indicating the absence of required parameters. In the example above, for instance, the application may display something like $\frac{\square}{\square}$ suggesting that a fraction was entered, but neither the numerator nor the denominator have been.

4.2.2 Pattern Mismatch

Rules $B.2$ and $B.3$ have been marked with a \dagger to indicate that the parser expects a token which is not found in the token stream. In both cases the parser will typically notify the user with a warning message.

4.2.3 Ambiguities

In TeX one cannot pass a macro with parameters as the parameter of another macro, unless the parameter is enclosed within a group. For example, it is an error to write `\sqrt{\sqrt{x}}`, the correct form is `\sqrt{\sqrt{x}}`. Because we treat the left curly brace like any other macro, grouping would not help our parser in resolving ambiguities. However, the

parser knows how many parameters a macro needs, because the token representing the control sequence has been annotated with such information by the lexer. When processing a macro with arguments the parser behaves “recursively”, it does not let an incomplete macro to be “captured” if it was passed as parameter of an outer macro. A consequence of this extension is that any well-formed fragment of \TeX markup is accepted by our parser resulting in the same structure, but there are some strings accepted by our parser that cause the \TeX parser to fail.

4.3 Incremental Parsing

Parsing must be efficient because it is performed in real-time, in principle at every modification of the input buffer, no matter how simple the modification is. Fortunately \TeX markup exhibits good *locality*, that is small modifications in the document cause small modifications in the parsing tree. Consequently we can avoid re-parsing the whole source document, we just need to re-parse a small interval of the input buffer around the point where the modification has occurred, and adjust the parsing tree accordingly. Let us consider again the example of Fig. 1 and suppose that a change is made in the markup

$$\{1\over{1+x}^2\} \Rightarrow \{1\over{1+x+y}^2\}$$

(a $+y$ is added to the denominator of the fraction). To be conservative we can re-parse the smallest term within braces that includes the modified part (the underlined fragments). Once the term has been re-parsed it has to be substituted in place of the old term in the parsing tree.

In order to compute the interval of the input buffer to be re-parsed we annotate the nodes of the parsing tree with information about the first and the last characters of the buffer which were scanned while building the node and all of its children. A simple visit of the tree can locate the smaller interval affected by the modification.

Curly braces occur frequently enough in the markup to give good granularity for re-parsing. At the same time limiting re-parsing to braced terms helps control the costs related to the visit to the parsing tree and to the implementation of the incremental parsing and transformation machinery.

5 Transformation

The transformation phase recognizes structured patterns in the parsing tree and generates corresponding fragments of the result document. We have already anticipated that XSLT is a very natural choice

for the implementation of this phase. Besides, XSLT stylesheets can be extended very easily, by providing new *templates* that recognize and properly handle new macros that an author has introduced.

We can see in Fig. 2 two sample templates taken from an XSLT stylesheet for converting the internal parsing tree into a MathML tree. Both templates have a preamble made of an `xs1:if` construct which we will discuss later in this section. Since the \TeX tree and the MathML tree are almost isomorphic (Fig. 1) the transformation is generally very simple and in many cases it amounts at just renaming the node labels. Template (a) is one such case: it matches any node in the parsing tree with label `macro` and having the `name` attribute set to `over`. The node for the `\over` macro corresponds naturally to the `mfrac` element in MathML. The two parameters of `\over` are transformed recursively by applying the stylesheet templates to the first and second child nodes (`p[1]` means “the first `p` child of this node”, similarly `p[2]` refers to the second `p` child).

Template (b) is slightly more complicated and shows one case where there is some change in the structure. For combined sub/super scripts \TeX accepts a sequence of `_` and `^` no matter in what order they occur, but MathML has a specific element for such expressions, namely `msubsup`. The template matches an `sb` node whose first parameter contains an `sp` node, thus detecting a `...^..._...` fragment of markup, then the corresponding `msubsup` element is created and its three children accessed in the proper position of the parsing tree. A symmetric template will handle the case where the subscript occurs before the superscript.

5.1 Incremental Transformation

As we have done for parsing, for transformations we also need to account for their cost. In a batch, one-shot conversion from \TeX this is not generally an issue, but in an interactive authoring tool a transformation is required at every modification of the parsing tree in order to update the view of the document.

Intuitively, we can reason that if only a fragment of the parsing tree has changed, we need re-transform only that fragment and substitute the result in the final document. This technique makes two assumptions: (1) that transformations are context-free; that is, the transformation of a fragment in the parsing tree is not affected by the context in which the fragment occurs; (2) that we are able to relate corresponding fragments between the parsing and the result trees.

<pre> <xsl:template match="macro[@name='over']"> <m:frac> <xsl:if test="@id"> <xsl:attribute name="xref"> <xsl:value-of select="@id"/> </xsl:attribute> </xsl:if> <xsl:apply-templates select="p[1]"/> <xsl:apply-templates select="p[2]"/> </m:frac> </xsl:template> </pre>	<pre> <xsl:template match="macro[@name='sb'] [p[1]/*[1][self::macro[@name='sp']]]"> <m:msubsup> <xsl:if test="@id"> <xsl:attribute name="xref"> <xsl:value-of select="@id"/> </xsl:attribute> </xsl:if> <xsl:apply-templates select="p[1]/*p[1]"/> <xsl:apply-templates select="p[2]"/> <xsl:apply-templates select="p[1]/*p[2]"/> </m:msubsup> </xsl:template> </pre>
(a)	(b)

Figure 2: Example of XSLT templates for the transformation of the internal parsing tree into a MathML tree. MathML elements can be distinguished because of the `m:` prefix.

Template (b) in Fig. 2 shows one case where the transformation is not context free: the deeper `sp` node is not processed as if it would occur alone, but it is “merged” together with its parent. More generally we can imagine that transformations can make almost arbitrary re-arrangements of the structure. This problem cannot be solved unless we make some assumptions, and the one we have already committed to in Sect. 4 is that braces define “black-box” fragments which can be transformed in isolation, without context dependencies.

As for the matter of relating corresponding fragments of the two documents, we use identifiers and references. Each node in the parsing tree is annotated with a unique identifier (in our sample templates we are assuming that the identifier is a string in the `id` attribute). Templates create corresponding `xref` attributes in the result document “pointing” to the fragment with the same identifier in the parsing tree. This way, whenever a fragment of the parsing tree is re-transformed, it replaces the fragment in the result document with the same identifier.

More generally, back-pointers provide a mechanism for relating the view of the document with the source markup. This way it is possible to perform operations like selection or cut-and-paste that, while having a visual effect in the view, act indirectly at the content/markup level.

6 Conclusion

We have presented architectural and implementation issues of an interactive editor based on \TeX syntax which allows flexible customization and content-oriented authoring. \TeX macs² is probably the existing application that most closely adopts such architecture, with the difference that \TeX macs does not

stick to \TeX syntax as closely as we do and that, apart from being a complete (and cumbersome) editing tool and not just an interface, it uses encoding and transformation technologies not based on standard languages (XML [2] and XSLT [3]).

Among batch conversion tools we observe a tendency to move towards the processing of content. The \TeX to MathML converter by Igor Rodionov and Stephen Watt at the University of Western Ontario [8, 9] is one such tool, and the recent Hermes converter by Romeo Anghelache [10] is another. These represent significant steps forwards when compared to converters such as \LaTeX 2HTML.³

A prototype tool called Edi \TeX , based on the architecture described in this paper, has been developed and is freely available along with its source code.⁴ No mention of MathML is made in the name of the tool to remark the fact that the architecture is very general and can be adapted to other kinds of markup. The prototype is currently being used as interface for a proof-assistant application where editing of complex mathematical formulas and proofs is required. In this respect we should remark that \TeX syntax is natural for “real” mathematics, but it quickly becomes clumsy when used for writing terms of programming languages or λ -calculus. This is mainly due to the conventions regarding spaces (for instance, spaces in the λ -calculus denote function application) and identifiers (the rule “one character is one identifier” is fine for mathematics, but not for many other languages). Note however that, since the lexical analyzer is completely separate from the rest of the architecture, the token stream being its interface, it can be easily targeted to a language with different conventions than those of \TeX .

² <http://www.texmacs.org/>

³ <http://www.latex2html.org/>

⁴ <http://helm.cs.unibo.it/software/editex/>

The idea of using some sort of restricted \TeX syntax for representing mathematical expressions is not new. For example, John Forkosh's \MimeTeX ⁵ generates bitmap images of expressions to be embedded in Web pages. However, to the best of our knowledge the formal specification of the parser for simplified \TeX markup presented in Sect. 4 is unique of its kind. A straightforward implementation based directly on the rules given in Table 2 amounts at only just 70 lines of functional code (in an ML dialect), which can be considered something of an achievement given that parsing \TeX is normally regarded as a hard task. By comparison, the parsing code in \MimeTeX amounts to nearly 350 lines of C code after stripping away the comments.

One may argue that the simplified \TeX markup is too restrictive, but in our view this is just the sensible fragment of \TeX syntax that the average user should be concerned about. In fact the remaining syntactic expressiveness provided by \TeX is mainly required for the implementation of complex macros and of system internals, which should never surface at the document level. By separating the transformation phase we shift the mechanics of macro expansion to a different level which can be approached with different (more appropriate) languages. Since this mode of operation makes the system more flexible we believe that our design is a valuable contribution which may provide an architecture for other implementers to adopt.

References

- [1] The Unicode Consortium: The Unicode Standard, Version 4.0, Boston, MA, Addison-Wesley (2003). <http://www.unicode.org/>
- [2] Tim Bray, Jean Paoli, C.M. Sperberg-McQueen, Eve Maler (editors): Extensible Markup Language (XML) 1.0 (2nd Edition), W3C Recommendation (2000). <http://www.w3.org/TR/2000/REC-xml-20001006>
- [3] James Clark (editor): XML Transformations (XSLT) Version 1.0, W3C Recommendation (1999). <http://www.w3.org/TR/1999/REC-xslt-19991116>
- [4] Ron Ausbrooks, Stephen Buswell, Stéphane Dalmas, Stan Devitt, Angel Diaz, et al.: Mathematical Markup Language (MathML) Version 2.0 (2nd Edition) W3C Recommendation, (2003). <http://www.w3.org/TR/2003/REC-MathML2-20031021/>
- [5] Donald E. Knuth: The \TeX book, Addison-Wesley, Reading, MA, USA (1994).
- [6] Leslie Lamport: A Document Preparation System: \LaTeX , Addison-Wesley, Reading, MA, USA (1986).
- [7] Douglas Lovell: \TeX XML: Typesetting XML with \TeX , *TUGboat*, 20(3), pp. 176–183 (September 1999).
- [8] Sandy Huerter, Igor Rodionov, Stephen M. Watt: Content-Faithful Transformations for MathML, Proc. International Conference on MathML and Math on the Web (MathML 2002), Chicago, USA (2002). <http://www.mathmlconference.org/2002/presentations/huerter/>
- [9] Stephen M. Watt: Conserving implicit mathematical semantics in conversion between \TeX and MathML, *TUGboat*, 23(1), pp. 108–108 (2002).
- [10] Romeo Anghelache: \LaTeX -based authoring tool, Deliverable D4.d, MoWGLI Project (2003). <http://relativity.livingreviews.org/Info/AboutLR/mowgli/index.html>

Appendix: The TML DTD

```

<!ENTITY % TML.node "
  empty|space|literal|macro">
<!ENTITY % TML.common.attrib "
  id          CDATA #IMPLIED
  xref        CDATA #IMPLIED
  start       NMTOKEN #IMPLIED
  end         NMTOKEN #IMPLIED">
<!ELEMENT empty EMPTY>
<!ATTLIST empty %TML.common.attrib;>
<!ELEMENT space EMPTY>
<!ATTLIST space
  %TML.common.attrib
  name      NMTOKEN #IMPLIED
  literal   CDATA #IMPLIED>
<!ELEMENT literal #PCDATA>
<!ATTLIST literal
  %TML.common.attrib;
  name      NMTOKEN #IMPLIED>
<!ELEMENT macro (p)*>
<!ATTLIST macro
  %TML.common.attrib;
  name      NMTOKEN #REQUIRED
  literal   CDATA #IMPLIED>
<!ELEMENT p (%TML.node;)*>
<!ATTLIST p %TML.common.attrib;>

```

⁵ <http://www.ctan.org/tex-archive/support/mimetex/>

Animations in pdfTeX-generated PDF

Jan Holeček

Faculty of Informatics, Masaryk University
Botanická 68a
60200 Brno
Czech Republic
holecek@fi.muni.cz
<http://www.fi.muni.cz/~xholecek>

Petr Sojka

Faculty of Informatics, Masaryk University
Botanická 68a
60200 Brno
Czech Republic
sojka@fi.muni.cz
<http://www.fi.muni.cz/usr/sojka>

Abstract

This paper presents a new approach for creating animations in Portable Document Format (PDF). The method of animation authoring described uses free software (pdfTeX) only. The animations are viewable by any viewer that supports at least some features of Acrobat JavaScript, particularly Adobe (Acrobat) Reader, which is available at no cost for a wide variety of platforms. Furthermore, the capabilities of PDF make it possible to have a single file with animations both for interactive viewing and printing.

The paper explains the principles of PDF, Acrobat JavaScript and pdfTeX needed to create animations for Adobe Reader using no other software except pdfTeX. We present a step by step explanation of animation preparation, together with sample code, using a literate programming style. Finally, we discuss other possibilities of embedding animations into documents using open standards (SVG) and free tools, and conclude with their strengths and weaknesses with respect to the method presented.

1 Introduction

Extensive use of electronic documents leads to new demands being made on their content. Developing specific document versions for different output devices is time consuming and costly. A very natural demand, especially when preparing educational materials, is embedding animations into a document.

A widely used open format for electronic documents is the Adobe PDF [2] format, which combines good typographic support with many interactive features. Even though it contains no programming language constructs such as those found in PostScript, the format allows for the inclusion of *Document Level JavaScript* (DLJS) [1]. Widely available PDF viewers such as Adobe Reader (formerly Acrobat Reader) benefit from this possibility, allowing interactive documents to be created.

One of the first applications showing the power of using JavaScript with PDF was Hans Hagen's calculator [5]. Further, the AcroTeX bundle [9] uses several L^AT_EX packages and the full version of the Adobe Acrobat software for preparing PDF files with DLJS [10]; macro support for animations is rudimentary and it is stressed in the documentation that it works only with the full commercial version of Acrobat.

Our motivation is a need for PDF animations in a textbook [3] published both on paper and on CD. We have published it using Acrobat [7, 8], and eventually discovered a method to create animations using pdfTeX [11] only.

pdfTeX facilitates the PDF creation process in several ways. We can directly write the PDF code which is actually required to insert an animation. We can also utilise the T_EX macro expansion power

to produce PDF code. And finally, we can write only the essential parts directly, leaving the rest to pdfTeX. pdfTeX introduces new primitives to take advantage of PDF features. The ones we are going to use will be described briefly as they appear.

In this paper, we present this new ‘pdfTeX only’ way of embedding animations. We require no previous knowledge either of the PDF language or of pdfTeX extensions to TeX. However, the basics of TeX macro definitions and JavaScript are assumed.

The structure of the paper is as follows. In the next section we start with the description of the PDF internal document structure with respect to animations. The core of the paper consists of commented code for the pdfTeX that generates a simple all-in-one animation. The examples are written in plain TeX [6], so that others can use it in elaborate macro packages, in a literate programming style. In the second example the animation is taken from an external file, allowing the modification of the animation without modifying the primary document. Finally, we compare this approach with the possibilities of other formats, including the new standard for Scalable Vector Graphics (SVG) [12] from the W3C.

2 The PDF Document Structure

A PDF file typically consists of a header, a body, a cross-reference table and a trailer. The body is the main part of the PDF document. The other parts provide meta-information and will not be discussed here. A PDF document is actually a graph of interconnected objects, each being of a certain type. There are basic data types (boolean, numeric, string) and some special and compound types which require some explanation.

A *name* object has the form /MYNAME. There is a set of names with predefined meanings when used as a dictionary key or value. Other names can be defined by the user as human readable references to indirect objects (dictionaries and indirect objects are treated below). An *array* object is a one-dimensional list, enclosed by square brackets, of objects not necessarily of the same type. A *dictionary* object is a hash, i.e., a set of key-value pairs where the keys are name objects and the values are arbitrary objects. A dictionary is enclosed by the << and >> delimiters. *Stream* objects are used to insert binary data into a PDF document. There is also a special *null* object used as an “undefined” value.

The body of a PDF file consists of a sequence of labelled objects called *indirect objects*. An object of any other type which is given a unique *object identifier* can form an indirect object. When an object is required in some place (an array element, a value

of a key in a dictionary), it can be given explicitly (a direct reference) or as an object identifier to an indirect object (an *indirect reference*). In this way objects are interconnected to form a graph. An indirect reference consists of two numbers. The first number is a unique object number. The second is an object version number and is always 0 in indirect objects newly created by pdfTeX — the first one therefore suffices to restore an indirect reference.

Various document elements are typically represented by dictionary objects. Each element has a given set of required and optional keys for its dictionary. For example, the document itself is represented by a *Catalog* dictionary, the root node of the graph. Its key-value pairs define the overall properties of the document. A brief description of concrete objects will be given when encountered for the first time. See [2] for more detailed information.

3 Insertion of the Animation Frames

We are not interested in constructing the animation frames themselves — any graphics program such as METAPOST will do. Let us hence assume we have a PDF file, each page of which forms a single animation frame and the frames are in the order of appearance.

Every image is inserted into PDF as a so-called *form XObject* which is actually an indirect stream object. There are three primitives that deal with images in pdfTeX. The `\pdfximage` creates an indirect object for a given image. The image can be specified as a page of another PDF file. However, the indirect object is actually inserted only if referred to by the `\pdfrefximage` primitive or preceded by `\immediate`. `\pdfrefximage` takes an object number (the first number of indirect reference) as its argument and adds the image to the TeX list being currently built. The object number of the image most recently inserted by `\pdfximage` is stored in the `\pdflastximage` register.

A general PDF indirect object can be created similarly by `\pdfobj`, `\pdfrefobj` and `\pdflastobj`. `\pdfobj` takes the object content as its argument. TeX macro expansion can be used for generating PDF code in an ordinary manner.

In our example, we first define four macros for efficiency. The `\ximage` macro creates a form XObject for a given animation frame (as an image) and saves its object number under a given key. The `\insertobj` macro creates a general PDF object and saves its object number under a given key. The `\oref` macro expands to an indirect reference of an object given by the argument. The last “R” is an operator that creates the actual indirect reference from

two numbers. We are not going to use `\pdfref*` primitives, so `\immediate` must be present. References will be put directly into the PDF code by the `\oref` macro. The `\image` macro actually places an image given by its key onto the page.

```

1  % an image for further use
2  \def\ximage#1#2{%
3    \immediate\pdfximage
4      page #2 {frames-in.pdf}%
5    \expandafter\edef
6      \csname pdf:#1\endcsname
7      {\the\pdfplastximage}}
8
9  % a general object for further use
10 \def\insertobj#1#2{%
11   \immediate\pdfobj{#2}%
12   \expandafter\edef
13     \csname pdf:#1\endcsname
14     {\the\pdflastobj}}
15
16 % expands to an indirect ref. for a key
17 \def\oref#1{%
18   \csname pdf:#1\endcsname\space 0 R}
19
20 % actually places an image
21 \def\image#1{%
22   \expandafter\pdfrefximage
23     \csname pdf:#1\endcsname}

```

Another new primitive introduced by pdfTeX is `\pdfcatalog`. Its argument is added to the document's Catalog dictionary every time it is expanded. The one below makes the document open at the first page and the viewer fit the page into the window. One more key will be described below.

```

24 % set up the document
25 \pdfcatalog{/OpenAction [ 0 /Fit ]}

```

Now we are going to insert animation frames into the document. We will use the `\ximage` macro defined above. Its first argument is the name to be bound with the resulting form XObject. The second one is the number of the frame (actually a page number in the PDF file with frames). One needs to be careful here because pdfTeX has one-based page numbering while PDF uses zero-based page numbering internally.

```

26 % all animation frames are inserted
27 \ximage{fr0}{1} \ximage{fr1}{2}
28 \ximage{fr2}{3} \ximage{fr3}{4}
29 \ximage{fr4}{5} \ximage{fr5}{6}
30 \ximage{fr6}{7} \ximage{fr7}{8}
31 \ximage{fr8}{9}

```

4 Setting up an AcroForm Dictionary

The interactive features are realized by annotation elements in PDF. These form a separate layer in addition to the regular document content. Each one denotes an area on the page to be interactive and binds some actions to various events that can happen for that area. Annotations are represented by Annot dictionaries. The way pdfTeX inserts annotations into PDF is discussed in the section “Animation Dynamics” below.

Annotations are transparent by default, i.e., the page appearance is left unchanged when adding an annotation. It is up to the regular content to provide the user with the information that some areas are interactive.

We will be interested in a subtype of annotations called *interactive form fields*. They are represented by a Widget subtype of the Annot dictionary. Widgets can be rendered on top of the regular content. However, some resources have to be set. The document's Catalog refers to an AcroForm dictionary in which this can be accomplished.

The next part of the example first defines the name `Helv` to represent the Helvetica base-font (built in font). This is not necessary but it allows us to have a smooth control button. Next we insert the AcroForm dictionary. The DR stands for “resource dictionary”. We only define the Font resource with one font. The DA stands for “default appearance” string. The `/Helv` sets the font, the `7 Tf` sets the font size scale factor to 7 and the `0 g` sets the color to be 0% white (i.e., black). The most important entry in the AcroForm dictionary is `NeedAppearances`. Setting it to `true` (line 43) makes the Widget annotations visible. Finally, we add the AcroForm dictionary to the document's Catalog.

```

32 % the Helvetica basefont object
33 \insertobj{Helv}{
34 << /Type /Font /Subtype /Type1
35   /Name /Helv
36   /BaseFont /Helvetica >> }
37
38 % the AcroForm dictionary
39 \insertobj{AcroForm}{
40 << /DR << /Font <<
41   /Helv \oref{Helv} >> >>
42   /DA (/Helv 7 Tf 0 g )
43   /NeedAppearances true >> }
44
45 % add a reference to the Catalog
46 \pdfcatalog{/AcroForm \oref{AcroForm}}

```

To make a form XObject with an animation frame accessible to JavaScript, it has to be assigned

a name. There are several namespaces in PDF in which this can be accomplished. The one searched for is determined from context. We are only interested in an AP namespace that maps names to annotation appearance streams. pdfTeX provides the `\pdfnames` primitive that behaves similarly to `\pdfcatalog`. Each time it is expanded it adds its argument to the Names dictionary referred from document's Catalog. The Names dictionary contains the name definitions for various namespaces. In our example we put definitions into a separate object `AppearanceNames`.

The name definitions may form a tree to make the lookup faster. Each node has to have Limits set to the lexically least and greatest names in its subtree. There is no extensive set of names in our example, so one node suffices. The names are defined in the array of pairs containing the name string and the indirect reference.

```

47 % defining names for frames
48 \insertobj{AppearanceNames}{
49 << /Names
50   [ (fr0) \oref{fr0} (fr1) \oref{fr1}
51     (fr2) \oref{fr2} (fr3) \oref{fr3}
52     (fr4) \oref{fr4} (fr5) \oref{fr5}
53     (fr6) \oref{fr6} (fr7) \oref{fr7}
54     (fr8) \oref{fr8} ]
55   /Limits [ (fr0) (fr8) ] >> }
56
57 % edit the Names dictionary
58 \pdfnames{/AP \oref{AppearanceNames}}

```

5 Animation Dynamics

We have created all the data structures needed for the animation in the previous section. Here we introduce the code to play the animation. It uses Acrobat JavaScript [1], an essential element of interactive forms. Acrobat JavaScript is an extension of Netscape JavaScript targeted to PDF and Adobe Acrobat. Most of its features are supported by Adobe Reader. They can, however, be supported by any other viewer. Nevertheless, the Reader is the only one known to us that supports interactive forms and JavaScript.

The animation is based on interchanging frames in a single widget. Here we define the number of frames and the interchange timespan in milliseconds to demonstrate macro expansion in JavaScript.

```

59 % animation properties
60 \def\frames{8}
61 \def\timespan{550}

```

Every document has its own instance of a JavaScript interpreter in the Reader. Every JavaScript

action is interpreted within this interpreter. This means that one action can set a variable to be used by another action triggered later. Document-level JavaScript code, e.g., function definitions and global variable declarations, can be placed into a JavaScript namespace. This code should be executed when opening the document.

Unfortunately, there is a bug in the Linux port of the Reader that renders this generally unusable. The document level JavaScript is not executed if the Reader is not running yet and the document is opened from a command line (e.g., `'acroread file.pdf'`). Neither the first page's nor the document's open action are executed, which means they cannot be used as a workaround. Binding a JavaScript code to another page's open action works well enough to suffice in most cases.

We redeclare everything each time an action is triggered so as to make the code as robust as possible. First we define the `Next` function, which takes a frame index from a global variable, increases it modulo the number of frames and shows the frame with the resulting index. The global variable is modified.

The animation actually starts at line 78 where the frame index is initialized. The frames are displayed on an interactive form's widget that we name `"animation"` — see "Placing the Animation" below. A reference to this widget's object is obtained at line 79. Finally, line 80 says that from now on, the `Next` function should be called every `\timespan` milliseconds.

```

62 % play the animation
63 \insertobj{actionPlay}{
64 << /S /JavaScript /JS (
65   function Next() {
66     g.delay = true;
67     if (cntr == \frames) {
68       cntr = 0;
69       try { app.clearInterval(arun); }
70         catch(except) {}
71     } else { cntr++; }
72     g.buttonsetIcon(
73       this.getIcon("fr" + cntr));
74     g.delay=false;
75   }
76   try { app.clearInterval(arun); }
77     catch(except) {}
78   var cntr = 0 ;
79   var g = this.getField("animation");
80   var arun = app.setInterval("Next()",
81                               \timespan);
82 ) >> }

```


Now, let us describe the `Next` function in more detail. Line 66 suspends widget’s redrawing until line 74. Then the global variable containing the current frame index is tested. If the index reaches the number of frames, it is set back to zero and the periodic calling of the function is interrupted. The function would be aborted on error, but because we catch exceptions this is avoided. The `getIcon` function takes a name as its argument and returns the reference to the appearance stream object according to the AP names dictionary. This explains our approach of binding the names to animation frames—here we use the names for retrieving them. The `buttonSetIcon` method sets the object’s appearance to the given icon.

Line 76 uses the same construct as line 69 to handle situations in which the action is relaunched even if the animation is not finished yet. It aborts the previous action. It would have been an error had the animation not been running, hence we must use the exception catching approach.

6 Placing the Animation

The animation is placed on an interactive form field—a special type of annotation. There are two primitives in pdfTeX, `\pdfstartlink` and `\pdfendlink`, to produce annotations. They are intended to insert hyperlink annotations but can be used for creating other annotations as well. The corresponding `\pdfstartlink` and `\pdfendlink` must reside at the same box nesting level. The resulting annotation is given the dimensions of the box that is enclosed by the primitives. We first create a box to contain the annotation. Note that both box and annotation size are determined by the frame itself—see line 91 where the basic frame is placed into the regular page content.

We will turn now to the respective entries in the annotation dictionary. The annotation is to be an interactive form field (`/Subtype /Widget`). There are many field types (FT). The only one that can take any appearance and change it is the *pushbutton*. It is a special kind of *button* field type (`/FT /Btn`). The type of button is given in an array of field bit flags `Ff`. The *pushbutton* has to have bit flag 17 set (`/Ff 65536`). To be able to address the field from JavaScript it has to be assigned a name. We have assigned the name `animation` to it as mentioned above (`/T (animation)`). Finally, we define the appearance characteristics dictionary `MK`. The only entry `/TP 1` sets the button’s appearance to consist only of an icon and no caption.

```
83 % an animation widget
84 \centerline{\hbox{%
```

```
85 \pdfstartlink user{
86 /Subtype /Widget /FT /Btn
87 /Ff 65536 /T (animation)
88 /BS << /W 0 >>
89 /MK << /TP 1 >> }%
90 \image{fr0}%
91 \pdfendlink}}
```

For the sake of brevity and clarity we are going to introduce only one control button in our example. However, we have defined a macro for creating control buttons to show a very simple way of including multiple control buttons. The `\controlbutton` macro takes one argument: the caption of the button it is to produce. The macro creates a *pushbutton* and binds it to an action defined like `actionPlay`.

We have chosen control buttons to be *pushbuttons* again. They are little different from the animation widget—they are supposed to look like buttons. The `BS` dictionary (i.e., border style) sets the border width to 1 point and style to 3D button look. The `MK` dictionary (appearance characteristics dictionary) sets the background color to 60% white and the caption (line 98). The `/H /P` entry tells the button to push down when clicked on. Finally, an action is bound to the button by setting the value of the `A` key.

```
92 % control button for a given action
93 \def\controlbutton#1{%
94 \hbox to 1cm{\pdfstartlink user{
95 /Subtype /Widget /FT /Btn
96 /Ff 65536 /T (Button#1)
97 /BS << /W 1 /S /B >>
98 /MK << /BG [0.6] /CA (#1) >>
99 /H /P /A \oref{action#1}
100 }\hfil\strut\pdfendlink}}
```

And finally, we add a control button that plays the animation just below the animation widget.

```
101 % control button
102 \centerline{\hfil
103 \controlbutton{Play}\hfil}
104
105 \bye
```

7 External Animation

Let us modify the example a little so that the animation frames will be taken from an external file. This has several consequences which will be discussed at the relevant points in the code.

We are going to completely detach the animation frames from the document. As a result, we will need only the `\insertobj` and `\oref` macros from lines 1–23 from the previous example. Lines 26–31 are no longer required.

A problem arises here: the basic frame should be displayed in the animation widget when the document is opened for the first time. This can be accomplished by modifying the `OpenAction` dictionary at line 25 as follows.

```
\pdfcatalog{ /OpenAction <<
  /S /JavaScript /JS (
    var g = this.getField("animation");
    g.buttonImportIcon(
      "frames-ex.pdf",0);
    this.pageNum = 0;
    this.zoomType = zoomtype.fitP;
  ) >> }
```

This solution suffers from the bug mentioned in the “Animation Dynamics” section. The animation widget will be empty until a user performs an action every time the bug comes into play.

We still do need an `AcroForm` dictionary, so lines 32–46 are left without a change. Lines 47–58 must be omitted on the other hand, as we have nothing to name. We are going to use the same animation as in the previous example, so lines 59–61 are left untouched. There is one modification of the JavaScript code to be done. The `buttonSetIcon` function call is to be replaced by

```
g.buttonImportIcon(
  "frames-ex.pdf", cntr);
```

We have used the basic frame to determine a size of the widget in the previous example. This is impossible now because it has to be done at compile time. The replacement for lines 83–91 is as follows

```
% an animation widget
\centerline{\hbox to 6cm{%
  \vrule height 6cm depth 0pt width 0pt
  \pdfstartlink user{
    /Subtype /Widget /FT /Btn
    /Ff 65536 /T (animation)
    /BS << /W 0 >>
    /MK << /TP 1
      /IF << /SW /A /S /P
        /A [0.5 0.5] >> >> }%
  \hfil\pdfendlink}}
```

Dimensions of the widget are specified explicitly and an `IF` (icon fit) dictionary is added to attributes of the pushbutton so that the frames would be always (`/SW /A`) proportionally (`/S /P`) scaled to fit the widget. Moreover, frames are to be centered in the widget (`/A [0.5 0.5]`) which would be the default behavior anyway. The basic frame is not placed into the document — there is only glue instead.

Lines 92–105 need not be modified.

8 Two Notes on Animation Frames

The examples with full \TeX source files can be found at <http://www.fi.muni.cz/~xholecek/animations/>. As one can see in these examples, the all-in-one approach allows all frames to share a single background which is formed by the frame actually inserted into the page. However, it is possible to overlay pushbuttons. Elaborate constructions, the simplest of which is to use a common background frame in the example with external animations, can be achieved in conjunction with transparency.

One must ensure the proper size of all frames when fitting them into the widget. We have encountered situations (the given example being one of them) where the bounding box of `METAPOST` generated graphics with \TeX label was not set properly using `\convertMPtoPDF` and a white line had to be drawn around the frames to force the proper bounding box as a workaround.

9 Animations in Other Formats

It is fair to list and compare other possible ways of creating animations. In this section we give a brief overview of a dozen other formats and technologies capable of handling animations.

9.1 GIF

One of the versions of the GIF format is the GIF89a format, which allows multi-image support, with bitmap only animations to be encoded within a single GIF file. GIF format supports transparency, interlacing and plain text blocks. It is widely supported in Internet browsers. However, there are licensing problems due to the compression methods used, and the format is not supported in freely available \TeX ware.

9.2 SWF

The SWF format by Macromedia allows storing frame-based animations, created e.g., by Macromedia’s Flash authoring tool. The SWF authoring tools have to compute all the animation frames at export time. As proprietary Flash plug-ins for a wide range of Internet browsers are available, animations in SWF are relatively portable. The power of SWF can be enriched with scripting by ActionScript. At the time of writing, we are not aware of any \TeX ware supporting SWF.

9.3 Java

One can certainly program animations in a general programming language like Sun’s Java. The drawback is that there are high demands on one’s

programming capabilities in Java when creating portable animations. With $\mathcal{N}\mathcal{T}\mathcal{S}$ (a \TeX reimplementation in Java), one can possibly combine \TeX documents with fully featured animations, at the expense of studying numerous available classes, interfaces and methods.

9.4 DOM

It is possible to reference every element in an HTML or XML document by means of the W3C's Document Object Model (DOM), a standard API for document structure.

DOM offers programmers the possibility of implementing animations with industry-standard languages such as Java, or scripting languages as ECMAScript, JavaScript or JScript.

9.5 SVG

The most promising language for powerful vector graphic animation description seems to be Scalable Vector Graphics (SVG), a W3C recommendation [12]. It is being developed for XML graphical applications, and since SVG version 1.1 there is rich support for animations. The reader is invited to look at the freely available book chapter [13] about SVG animations on the publisher's web site, or reading [4] about the first steps of SVG integration into \TeX world. There are freely available SVG viewers from Adobe (browser plug-in), Corel, and the Apache Foundation (Squiggle).

SVG offers even smaller file sizes than SWF or our method. The description of animations is time-based, using another W3C standard, SMIL, Synchronised Multimedia Integration Language. The author can change only one object or its attribute in the scene at a time, allowing detailed control of animated objects through the declarative XML manner. Compared to our approach, this means a much wider range of possibilities for creators of animations.

The SVG format is starting to be supported in \TeX ware. There are SVG backends in V \TeX and BaKoMa \TeX , and a program Dvi2Svg by Adrian Frischauf, available at <http://www.activemath.org/~adrianf/dvi2svg/>. Another implementation of a DVI to SVG converter in C is currently being developed by Rudolf Sabo at the Faculty of Informatics, Masaryk University in Brno.

10 Conclusions

We have shown a method of preparing both space-efficient and high-quality vector frame-based animations in PDF format using only freely available, \TeX -integrated tools.

11 Acknowledgments

Authors thank Oleg Alexandrov and Karl Berry for comments on an early draft of the paper.

The work has been supported by VZ MSM 143300003.

References

- [1] Adobe Systems Incorporated. Acrobat JavaScript Object Specification, Version 5.1, Technical Note #5186. Technical report, Adobe, 2003. <http://partners.adobe.com/asn/developer/pdfs/tn/5186AcroJS.pdf>.
- [2] Adobe Systems Incorporated. *PDF Reference: Adobe Portable Document Format Version 1.5*. Addison-Wesley, Reading, MA, USA, fourth edition, August 2003.
- [3] Zuzana Došlá, Roman Plch, and Petr Sojka. Mathematical Analysis with Maple: 2. Infinite Series. CD-ROM, <http://www.math.muni.cz/~plch/nkpm/>, December 2002.
- [4] Michel Goossens and Vesa Sivunen. \LaTeX , SVG, Fonts. *TUGboat*, 22(4):269–280, October 2001.
- [5] Hans Hagen. The Calculator Demo, Integrating \TeX , METAPOST, JavaScript and PDF. *TUGboat*, 19(3):304–310, September 1998.
- [6] Petr Olšák. *\TeX book naruby (in Czech)*. Konvoj, Brno, 1997.
- [7] Petr Sojka. Animations in PDF. In *Proceedings of the 8th Annual Conference on Innovation and Technology in Computer Science Education, ITiCSE 2003*, page 263, Thessaloniki, 2003. Association of Computing Machinery.
- [8] Petr Sojka. Interactive Teaching Materials in PDF using JavaScript. In *Proceedings of the 8th Annual Conference on Innovation and Technology in Computer Science Education, ITiCSE 2003*, page 275, Thessaloniki, 2003. Association of Computing Machinery.
- [9] Donald P. Story. Acro \TeX : Acrobat and \TeX team up. *TUGboat*, 20(3):196–201, Sep. 1999.
- [10] Donald P. Story. Techniques of introducing document-level JavaScript into a PDF file from \LaTeX source. *TUGboat*, 22(3):161–167, September 2001.
- [11] Hán Th e Th anh. Micro-typographic extensions to the \TeX typesetting system. *TUGboat*, 21(4):317–434, December 2000.
- [12] W3C. *Scalable Vector Graphics (SVG) 1.1 Specification*, January 2003.
- [13] Andrew H. Watt. *Designing SVG Web Graphics*. New Riders Publishing, September 2001.

Arabic Mathematical e-Documents

Mustapha Eddahibi
m.eddahibi@ucam.ac.ma

Azzeddine Lazrek
lazrek@ucam.ac.ma

Khalid Sami
k.sami@ucam.ac.ma
Department of Computer Sciences,
Faculty of Sciences,
University Cadi Ayyad
P.O. Box 2390,
Marrakech, Morocco
<http://www.ucam.ac.ma/fssm/RyDARab>

Abstract

What problems do e-documents with mathematical expressions in an Arabic presentation present? In addition to the known difficulties of handling mathematical expressions based on Latin script on the web, Arabic mathematical expressions flow from *right to left* and use *specific symbols* with a *dynamic cursivity*. How might we extend the capabilities of tools such as MathML in order to structure Arabic mathematical e-documents? Those are the questions this paper will deal with. It gives a brief description of some steps toward an extension of MathML to mathematics in Arabic exposition. In order to evaluate it, this extension has been implemented in Mozilla.

KEYWORDS: Mathematical expressions, Arabic mathematical presentation, Multilingual documents, e-documents, Unicode, MathML, Mozilla.

1 Overview

It is well known that HTML authoring capabilities are limited. For instance, mathematics is difficult to search and web formatting is poor. For years, most mathematics on the web consisted of texts with scientific notation rendered as images. Image-based equations are generally harder to see, read and comprehend than the surrounding text in the browser window. Moreover, the large size of this kind of e-document can represent a serious problem. These problems become worse when the document is printed. For instance, the resolution of the equations will be around 72 dots per inch, while the surrounding text will typically be 300 or more dots per inch. In addition to the display problems, there are encoding difficulties. Mathematical objects can neither be searched nor exchanged between software systems nor cut and pasted for use in different contexts nor verified as being mathematically correct. As mathematical e-documents may have to

be converted to and from other mathematical formats, they need encoding with respect to both the mathematical notation and mathematical meaning.

The mathematical markup language MathML [14] offers good solutions to the previous problems. MathML is an XML application for describing mathematical notation and capturing both its structure, for high-quality visual display, and content, for more semantic applications like scientific software. XML stands for eXtensible Markup Language. It is designed as a simplified version of the meta-language SGML used, for example, to define the grammar and syntax of HTML. One of the goals of XML is to be suitable for use on the web by separating the presentation from the content. At the same time, XML grammar and syntax rules carefully enforce document structure to facilitate automatic processing and maintenance of large document collections.

MathML enables mathematics to be served, received, and processed on the web, just as HTML

has enabled this functionality for text. MathML elements can be included in XHTML documents with namespaces and links can be associated to any mathematical expression through XLink. Of course, there are complementary tools. For instance, the project OpenMath [12] also aims at encoding the semantics of mathematics without being in competition with MathML.

Now, what about some of the MathML internationalization aspects—say, for instance, its ability to structure and produce e-documents based on non Latin alphabets, such as mathematical documents in Arabic?

2 Arabic Mathematical Presentation

Arabic script is cursive. Small curves and ligatures join adjacent letters in a word. The shapes of most of the letters are context dependent; that is, they change according to their position in the word. Certain letters have up to four different shapes.

Although some mathematical documents using Arabic-based writing display mathematics in Latin characters, in general, not only the text is encoded with the Arabic script but mathematical objects and expressions are also encoded with special symbols flowing from right to left according to the Arabic writing. Moreover, some of these symbols are extensible.

Mathematical expressions are for the most part handwritten and introduced as images. A highly-evolved system of calligraphic rules governs Arabic handwriting. Though Arabic mathematical documents written by hand are sometimes of fair quality, the mere presentation of scientific documents is no longer enough, since there is a need for searchability, using them in software and so on.

The RyDArab [8] system makes it possible to compose Arabic mathematical expressions of high typographical quality. RyDArab complements \TeX for typesetting Arabic mathematical documents. RyDArab uses the Computer Modern fonts and those of Ω [4] or Arab \TeX [7]. The output is DVI, PS, PDF or HTML with mathematical expressions as images. The RyDArab [2] system does not replace or modify the functionality of the \TeX engine, so it does not restrict in any way the set of macros used for authoring. Automatic translation from and to Latin-based expressions is provided beginning with the latest RyDArab version. Will this be enough to structure and typeset e-documents with mathematics even when they are based on an alternative script? Starting from this material with \TeX and Ω , will MathML be able to handle Arabic mathematics?

3 MathML and Arabic Mathematics

Of course, semantically speaking, an Arabic mathematical expression is the same as a Latin-based one. Thus, only display problems need be taken into account. In any way, encoding semantics are beyond the scope of this paper.

In order to know if there really is a need to construct a new tool or only to improve an already available one, what are the possibilities offered by the known MathML renderers? As much of the work is built around \TeX , an *open source community effort*, it is hard to be precise about the current status of all \TeX /MathML related projects. Most of these projects belong to one of three basic categories:

- Conversions from \TeX to MathML. Of particular note here, are Ω [5, 6] and $\TeX4ht$ [13], a highly specialized editor/DVI driver. Both of these systems are capable of writing presentation MathML from \TeX documents. There are other converters such as $\LaTeX2HTML$ and $tralics$ [1].
- Conversions from MathML to \TeX . The conversion from MathML to \TeX can be done for instance, through reading MathML into Mathematica or other similar tools and then saving the result back out as \TeX , or using Scientific WorkPlace for suitable \LaTeX sources. The Con \TeX t system is another example.
- Direct typesetting of MathML using \TeX .

Currently, MathML is supported by many applications. This fact shows not only that it is the format of choice for publishing equations on the web but also that it is a universal interchange format for mathematics. More than twenty implementations are listed on the MathML official website, showing that all categories of mathematical software can handle MathML. Actually,

- most mathematical software, such as Scientific WorkPlace, Maple, MathCad and Mathematica, can export and import MathML;
- all common browsers can display MathML equations either natively or through the use of plug-ins;
- editors such as MathType, Amaya, \TeX macs, and WebEQ support MathML.

Once non-free or non-open-source tools are omitted, two web browsers remain: the well-known Mozilla system [11] and Amaya. The W3C's Amaya editor/browser allows authors to include mathematical expressions in web pages, following the MathML specification. Mathematical expressions are handled as structured components, in the same way

and in the same environment as HTML elements. All editing commands provided by Amaya for handling text are also available for mathematics, and there are some additional controls to enter and edit mathematical constructs. Amaya shows how other W3C specifications can be used in conjunction with MathML.

In the end, we chose to adapt Mozilla to the needs of the situation, mainly because of its popularity and widespread adoption as well as the existence of an Arabic version. The layout of mathematical expressions in Latin writing, and consequently that of the mathematical documents in Mozilla is more elegant and of good typographical quality compared to other systems.

For this implementation, we used the Mozilla 1.5 C++ source under Linux. Until now, there was no Mozilla version with support for bidirectionality or cursivity in a mathematical environment. In math mode, only left to right arrangement is supported. Thus, the first step is to find out how to get bidirectionality and cursivity inside a MathML passage.

In fact, adding the property of bidirectionality to MathML elements is a delicate task. It requires a careful study of the various possible conflicts. The bidirectionality algorithm for mathematical expressions is probably different from that originally in use for text.

Now, let us have a look at what would happen if the bidirectionality algorithm for HTML were used for MathML elements.

The MathML expression

```
<mn>1</mn>
<mo>+</mo>
<mi>ب</mi>
<mo>-</mo>
<mn>2</mn>
```

will be rendered as $1+2-ب$ instead of the expected equation: $1+ب-2$.

Since XML supports Unicode, we might expect that the introduction of Arabic text into MathML encoding would go without any problem. In other words, the Arabic text would be rendered from right to left, and letters would be connected just as they should be in their cursive writing. Will the use of the element `<mtext>` (similar to the use of the \TeX command `\hbox`) be enough to get a satisfactory rendering of Arabic?

The following Arabic text is a sample of what is obtained with `<mtext>` in Mozilla:

```
<mtext>
نص رياضي
</mtext>
```

رياضي صدين

The following Arabic abbreviation of the cosine function is an example of what we get if we introduce it with `<mi>`:

```
<mi>جتا</mi>
```

اتج

In order to allow the arrangement of sub-expressions from right to left in a given mother expression, a new element denoted `<rl>` is introduced.¹

```
<mrow>
<rl>
<mi>ب</mi>
<mo>+</mo>
<mi>س</mi>
</rl>
</mrow>
```

ب + س

The use of the element `<rl>` also allows solving the previous problem of introducing Arabic text in a mathematical expression.

```
<mtext>
<rl>نص رياضي</rl>
</mtext>
```

نص رياضي

```
<mi><rl>جتا</rl></mi>
```

جتا

The element `<rl>` can be used to transform some mathematical objects, such as left/right or open/close parentheses, into their mirror image.

```
<rl>
<mo><rl>[</rl></mo>
<mi>ب</mi>
<mo>,</mo>
<mn>3</mn>
<mo><rl>]</rl></mo>
</rl>
```

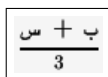
(3 , ب]

We can remark here that the symbol “,” has not changed to its mirror image “;”. The result is the same even when the comma is governed by `<rl>` (i.e., `<rl>`, `</rl>`). This symbol is not yet mentioned in the Bidi Mirroring list in the Unicode Character Database.

Particular arrangement of the arguments is made necessary by the MathML renderer for any presentation elements requiring more than one argument. On the other hand, elements of vertical arrangement such as `<mfrac>` do not need special handling.

¹ The name `rl` reminds us of the initials of right-to-left. Furthermore, because of the expected heavy use of this element, its name should be as short as possible.

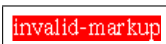
```
<mfrac>
<rl>
<mi>ب</mi>
<mo>+</mo>
<mi>س</mi>
</rl>
<mn>3</mn>
</mfrac>
```



Although the addition of `<rl>` helps to get Arabic text rendered as expected and to solve arrangement of sub-expressions within an expression, for certain elements, it doesn't work.

Using the element `<rl>` to get a superscript element `<msup>` or a subscript element `<msub>`, in the suitable right to left positions, generates a syntax error because `<msup>` requires two arguments, whereas there is only one argument, as can be seen in the following example:

```
<msup>
<rl>
<mi>س</mi>
<mi>ب</mi>
</rl>
</msup>
```



In this case, we introduce a new markup element `<amsup>`.² It changes the direction of rendering expressions while keeping the size of superscripts as it is with `<msup>`.

```
<amsup>
<mi>س</mi>
<mi>ب</mi>
</amsup>
```



The same principle is applied to other elements like `<msub>`. The notation of the arrangement in the Arabic combination analysis is different from its Latin equivalent.

```
<amarrange>
<mi>ل</mi>
<mn>5</mn>
<mn>2</mn>
</amarrange>
```

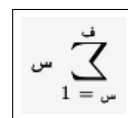


The next step is related to the shape of some symbols. In Arabic mathematical presentation, certain symbols, such as the square root symbol or the sum, in some Arabic areas, are built through a symmetric reflection of the corresponding Latin ones. These symbols require first the introduction of a new font family such as the one offered in the Arabic

² The following new elements defined in this system are prefixed with the initial of Arabic "a".

Computer Modern fonts. This family corresponds to the Computer Modern fonts with a mirror effect on some glyphs. In the same way that the Computer Modern fonts are used in the Latin mathematical environment with `<math>` the new element `<amath>` will allow the use of the Arabic Computer Modern fonts in the Arabic mathematical environment. The element `<amath>` would not be necessary if the Arabic mathematical symbols were already added in the Unicode tables. In fact, we use the same entity name and code for some symbols and their mirror images used in the Arabic presentation. For example, the Unicode name N-ARY SUMMATION coded by U+02211 is associated simultaneously to the Latin sum³ symbol \sum and to its Arabic equivalent mirror \sum . Thus, to specify which glyph, and consequently which font, is called, the introduction of a new element `<amath>` is necessary. This element would not be necessary if the symbols were denoted with two different entity names and consequently two different codes.

```
<amath>
<rl>
<mstyle displaystyle="true">
<munderover>
<mo>&sum;</mo>
<mrow>
<rl>
<mi>س</mi>
<mo>=</mo>
<mn>1</mn>
</rl>
</mrow>
<mi>ف</mi>
</munderover>
</mstyle>
<mi>س</mi>
</rl>
</amath>
```



In order to distinguish alphabetical symbols, in different shapes, from letters used in Arabic texts, and to avoid the heterogeneity resulting from the use of several fonts, there is a need for a complete Arabic mathematical font. That's exactly what we are trying to do in another project discussed elsewhere in this volume [10]. While waiting for their adoption by Unicode, the symbols in use in this font will be located in the Private Use Area E000-F8FF in the Basic Multilingual Plane.

```
<mi>&#xE004;</mi>
```



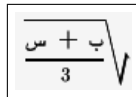
The use of the Arabic Computer Modern fonts is not enough for composed symbols. For example,

³ introduced as `∑` or `∑`.

the square root symbol is composed of the square root glyph supplemented by an over bar. This over bar is added by the renderer, which, thanks to a calculation of the width of the base, gives the length of this over bar.

In this case, neither the inversion of the glyph nor the use of the right-to-left element `<rl>` changes the direction of the visual rendering of the square root. For this reason we have introduced a new element (`<amsqrt>`), which uses the square root glyph from the Arabic Computer Modern font that shows the over bar to its left.

```
<amath>
  <amsqrt>
    <mfrac>
      <rl>
        <mi>ب</mi>
        <mo>+</mo>
        <mi>س</mi>
      </rl>
      <mn>3</mn>
    </mfrac>
  </amsqrt>
</amath>
```



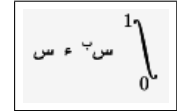
The root element `<amroot>` requires a treatment similar to that of the square root combined with the positioning of the index on the right of the baseline.

```
<amath>
  <amroot>
    <mi>س</mi>
    <mn>3</mn>
  </amroot>
</amath>
```



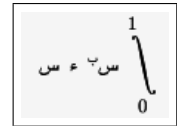
For the elements `<munderover>`, `<munder>`, `<mover>`, and `<mssubsup>`, italic correction needs to be done. In fact, mathematical symbols like the integral are slanted and the indices and exponents are shifted in the direction of the symbol's slant. This fact appears clearly in the following example representing two integrals, while using `<amsubsup>` in the first:

```
<amath>
  <rl>
    <mstyle displaystyle="true">
      <amsubsup>
        <mo>&int;</mo>
        <mn>0</mn>
        <mn>1</mn>
      </amsubsup>
    </mstyle>
    <amsup>
      <mi>س</mi>
    </amsup>
    <mi>ب</mi>
  </rl>
  <amath>
    <mi>س</mi>
    <mi>ء</mi>
    <mi>س</mi>
  </amath>
```

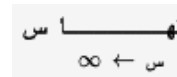


or `<amunderover>` in the second:

```
<amath>
  <rl>
    <mstyle displaystyle="true">
      <amunderover>
        <mo>&int;</mo>
        <mn>0</mn>
        <mn>1</mn>
      </amunderover>
    </mstyle>
    <amsup>
      <mi>س</mi>
    </amsup>
    <mi>ب</mi>
  </rl>
  <amath>
    <mi>س</mi>
    <mi>ء</mi>
    <mi>س</mi>
  </amath>
```

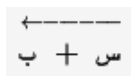


For the limit of an expression, manual lengthening of the limit symbol is performed. Of course, dynamic lengthening via automatic calculation of the width of the text under the limit sign would be better.



A lengthening of the straight line is not in conformity with the rules of the Arabic typography. A curvilinear lengthening is required, which can be obtained by using CurExt [9], which makes it possible to stretch Arabic letters according to calligraphic rules.

The following mathematical expression is an example of the use of `<mover>`, with automatic lengthening of the over arrow.



In fact, the use of the element `<r1>` doesn't represent a very practical solution as the encoding becomes heavier. The addition of this element must be transparent for the user; the same for all other new elements since they affect only the presentation and not the semantics of expression. An alternative solution consists of either building a new algorithm of bidirectionality for mathematics, or of adding attributes that will make it possible to choose the mathematical notation of the expression. We intend to use a new attribute `nota` for the root element `<math>`. It would indicate whether Arabic or Latin is used inside the mathematical expression. As the layout of a mathematical expression follows a precise logic, the direction of writing would be handled automatically without requiring the use of direction attributes for each child of the element `<math>`.

The FIGUE [3] system is an engine for the interactive rendering of structured objects. It allows the rendering of an Arabic text from right to left including some Latin mathematical expressions flowing from left to right thanks to a proposed bidirectional extension of MathML.

4 Conclusion

Our goal was to identify the difficulties and limitations that might obstruct the use of MathML for writing mathematics in Arabic. The main adaptation we made to MathML for Arabic mathematics was the addition of the element `<r1>` that allows:

- writing mathematical expressions from right-to-left;
- the use of specific symbols thanks to the modification of other elements;
- and handling the cursivity of writing.

Now, Arabic mathematical e-documents can be structured and published on the web using this extended version of Mozilla. Such documents can thus benefit from all the advantages of using MathML. Our project for the development of communication and publication tools for scientific and technical e-documents in Arabic is still at its beginning. We hope that the proposals contained in this paper will help to find suitable recommendations for Arabic mathematics in Unicode and MathML.

References

- [1] <http://www-sop.inria.fr/miaou/tralics/>.
- [2] Mustapha Eddahibi and Azzeddine Lazrek, *Arabic scientific document composition*, International Conference on Information Technology and Natural Sciences (ICITNS 2003, Amman, Jordan), 2003.
- [3] Hanane Naciri et Laurence Rideau, *Affichage et diffusion sur Internet d'expressions en langue arabe de preuves mathématiques*, CARI 2002 (Cameroun), 2002.
- [4] Yannis Haralambous and John Plaice, Multilingual Typesetting with Ω , a Case Study: Arabic, *Proceedings of the International Symposium on Multilingual Information Processing (Tsukuba)*, 1997, pp. 137–154.
- [5] Yannis Haralambous and John Plaice, Produire du MathML et autres *ML à partir d' Ω : Ω se généralise, *Cahiers GUTenberg*, vol. 33-34, 1999, pp. 173–182.
- [6] Yannis Haralambous and John Plaice, *X_{La}TeX, a DTD/Schema Which is Very Close to L_ATeX*, EuroTeX 2003: 14th European TeX Conference (ENST Bretagne, France), 2003 (to appear in *TUGboat*).
- [7] Klaus Lagally, *ArabTeX — Typesetting Arabic with Vowels and Ligatures*, EuroTeX'92 (Prague), 1992.
- [8] Azzeddine Lazrek, A package for typesetting arabic mathematical formulas, *Die TeXnische Komödie*, DANTE e.V., vol. 13. (2/2001), 2001, pp. 54–66.
- [9] Azzeddine Lazrek, *CurExt, Typesetting variable-sized curved symbols*, EuroTeX 2003 preprints: 14th European TeX Conference (Brest, France), 2003, pp. 47–71 (to appear in *TUGboat*).
- [10] Mostafa Banouni, Mohamed Elyaakoubi and Azzeddine Lazrek, *Dynamic Arabic mathematical fonts*, International Conference on TeX, XML and Digital Typography (TUG 2004, Xanthi, Greece), 2004.
- [11] Mozilla, <http://www.mozilla.org>.
- [12] OpenMath, <http://www.openmath.org/>.
- [13] TeX4ht, <http://www.cis.ohio-state.edu/~gurari/TeX4ht/mn.html>.
- [14] Presentation MathML and Content MathML, <http://www.w3.org/TR/MathML2>.

Dynamic Arabic Mathematical Fonts

Mostafa Banouni
mbanouni@voila.fr

Mohamed Elyaakoubi
m.elyaakoubi@ucam.ac.ma

Azzeddine Lazrek
lazrek@ucam.ac.ma
Department of Computer Sciences
Faculty of Sciences
University Cadi Ayyad
P.O. Box 2390
Marrakech, Morocco
<http://www.ucam.ac.ma/fssm/rydarab>

Abstract

This contribution describes a font family designed to meet the requirements of typesetting mathematical documents in an *Arabic presentation*. Thus, not only is the text written in an Arabic alphabet-based script, but specific symbols are used and mathematical expressions also spread out from right to left. Actually, this font family consists of two components: an Arabic mathematical font and a dynamic font. The construction of this font family is a first step of a project aiming at providing a complete and homogeneous Arabic font family, in the OpenType format, respecting Arabic calligraphy rules.

KEYWORDS: Mathematical font, Dynamic font, Variable-sized symbols, Arabic mathematical writing, Multilingual documents, Unicode, POSTSCRIPT, OpenType.

1 Overview

The Arabic language is native for roughly three hundred million people living in the Middle East and North Africa. Moreover, the Arabic script is used, in various slightly extended versions, to write many major languages such as Urdu (Pakistan), Persian and Farsi (Iran, India), or other languages such as Berber (North Africa), Sindhi (India), Uyghur, Kirgiz (Central Asia), Pashtun (Afghanistan), Kurdish, Jawi, Baluchi, and several African languages. A great many Arabic mathematical documents are still written by hand. Millions of learners are concerned in their daily learning by the availability of systems for typesetting and structuring mathematics.

Creating an Arabic font that follows calligraphic rules is a complex artistic and technical task, due in no small part to the necessity of complex contextual analysis. Arabic letters vary their form according to their position in the word and according to the neighboring letters. Vowels and diacrit-

ics take their place over or under the characters, and that is also context dependent. Moreover, the *kashida*, a small flowing curve placed between Arabic characters, is to be produced and combined with characters and symbols. The *kashida* is also used for the text justification. The techniques for managing the *kashida* are similar to those that can be used for drawing curvilinear extensible mathematical symbols, such as *sum*, *product* or *limit*.

There are several Arabic font styles. Of course, it is not easy to make available all existing styles. The font style *Naskh* was the first font style adopted for computerization and standardization of Arabic typography. So far, only Naskh, Koufi, Ruqaa, and to a limited extent Farsi have really been adapted to the computer environment. Styles like Diwani or Thuluth, for example, don't allow enough simplification, they have a great variation in characters shapes, the characters don't share the same baseline, and so on. Considering all that, we have decided to use the Naskh style for our mathematical font.

The RyDArab [10] system was developed for the purpose of typesetting Arabic mathematical expressions, written from right to left, using specific symbols. RyDArab is an extension of the $\text{T}_{\text{E}}\text{X}$ system. It runs with K. Lagally's Arabic system Arab $\text{T}_{\text{E}}\text{X}$ [8] or with Y. Haralambous and J. Plaiçe's multilingual Ω [6] system. The RyDArab system uses characters belonging to the Arab $\text{T}_{\text{E}}\text{X}$ font `xnsh` or to the `omsea` font of Ω , respectively. Further Arabic alphabetic symbols in different shapes can be brought from the font `NasX` that has been developed, for this special purpose, using METAFONT. The RyDArab system also uses symbols from Knuth's Computer Modern family, obtained through adaptation to the right-to-left direction of Arabic.

Since different fonts are in use, it is natural that some heterogeneity will appear in mathematical expressions typeset with RyDArab [9]. Symbol sizes, shapes, levels of boldness, positions on the baseline will not quite be in harmony. So, we undertook building a new font in OpenType format with two main design goals: on the one hand, all the symbols will be drawn with harmonious dimensions, proportions, boldness, etc., and on the other hand, the font should contain the majority of the symbols in use in the scientific and technical writing based on an Arabic script.

Both Arabic texts and mathematical expressions need some additional variable-sized symbols. We used the CurExt [11] system to generate such symbols. This application was designed to generate automatically curvilinear extensible symbols for $\text{T}_{\text{E}}\text{X}$ with the font generator METAFONT. The new extension of CurExt does the same with the font generator POSTSCRIPT.

While METAFONT generates bitmap fonts and thus remains inside the $\text{T}_{\text{E}}\text{X}$ environment, OpenType [14] gives outline and multi-platform fonts. Moreover, since Adobe and Microsoft have developed it jointly, OpenType has become a standard combining the two technologies TrueType and POSTSCRIPT. In addition, it offers some additional typographic layout possibilities thanks to its multi-table feature.

2 A Mathematical Font

The design and the implementation of a *mathematical font* are not easy [5]. It becomes harder when it is oriented to Arabic presentation. Nevertheless, independent attempts to build an Arabic mathematical font have been undertaken. In fact, F. Alhargan [1] has sent us proofs of some Arabic mathematical symbols in TrueType format.

Now we will describe the way we constructed the OpenType Arabic mathematical font `RamzArab`. The construction of the font started by drawing the whole family of characters by hand. This task was performed by a calligrapher. Then the proofs were scanned to transform them into vectors. The scanning tools alone don't produce a satisfying result, so once the design is finalized, the characters are processed and analyzed using special software to generate the file defining the font.

In Arabic calligraphy, the feather's head (*kalam*) is a flat rectangle. The writer holds it so that the largest side makes an angle of approximately 70° with the baseline. Except for some variations, this orientation is kept all along the process of drawing the character. Furthermore, as Arabic writing goes from right to left, some boldness is produced around segments from top left toward the bottom right and conversely, segments from top right to the bottom left will rather be slim as in Figure 1.

The `RamzArab` font in Figure 4 contains only symbols specific to Arabic mathematics presentation plus some usual symbols found even in text mode. It is mainly composed of the following symbols:

- alphabetic symbols: Arabic letters in various forms, such as characters in isolated standard form, isolated double-struck, initial standard, initial with tail, initial stretched and with loop (e.g., ب ب ب ب ب ب respectively);
- punctuation marks (e.g., ، ؛ : ! ؟);
- digits as used in the Maghreb Arab (North Africa), and as they are in the Machreq Arab (Middle East);
- accents to be combined with alphabetic symbols (e.g., ◌ ◌ ◌);
- ordinary mathematical symbols such as delimiters, arithmetic operators, etc.
- mirror image of some symbols such as sum, integral, etc.

In Arabic mathematics, the order of the alphabetic symbols differs from the Arabic alphabetic order. Some problems can appear with the alphabetic symbols in their multi-form.

Generally, in Arabic mathematical expressions, alphabetic symbols are written without dots (e.g., ◌ ◌ ◌) or diacritics. This helps to avoid confusions with accents. The dots can be added whenever they are needed, however. Thus, few symbols are left.

Moreover, some deviation from the general rules will be necessary: in a mathematical expression, the isolated form of the letter ALEF can be confused with the Machreq Arab digit ONE. The isolated

form of the letter HEH can also present confusion with the Machreq Arab digit FIVE. The choice of the glyphs ه and ه to denote respectively these two characters will help to avoid such confusions. Even though these glyphs are not in conformity with the homogeneity of the font style and calligraphic rules, they are widely used in mathematics. In the same way, the isolated form of the letter KAF ك , resulting from the combination of two other basic elements, will be replaced by the KAF glyph in Ruqaa style, ك .

For the four letters ALEF, DAL, REH and WAW, the initial and the isolated forms are the same, and these letters will be withdrawn from the list of letters in initial form. On the other hand, instead of a unique cursive element, the stretched form of each of the previous letters will result from the combination of two elements. It follows that these letters will not be present in the list of the letters in the stretched form.

The stretched form of a letter is obtained by the addition of a MADDA-FATHA or ALEF in its final form ل to the initial form of the letter to be stretched (e.g., $\text{د} + \text{ل} \rightarrow \text{دا}$). The glyph of LAM-ALEF لا has a particular ligature that will be added to the list. The stretched form of a character is used if there is no confusion with any usual function abbreviation (e.g., ح or جا for the sine function).

The form with tail is obtained starting from the initial form of the letter followed by an alternative of the final form of the letter HEH \sim (e.g., $\text{د} + \sim \rightarrow \text{د}\sim$). These two forms are not integrated into the font because they can be obtained through a simple composition.

The form with loop is another form of letters with a tail. It is obtained through the combination of the final form with a particular curl that differs from one letter to another (e.g., صه مَع). This form will be integrated into the font because it cannot be obtained through a simple composition.

The following particular glyphs are also in use:
 $\text{لا ء م ك ل ع ر ه ه ه}$

The elements that are used in the composition of the operator sum, product, limit and factorial in a conventional presentation (لها حد مح) are added also. These symbols are extensible. They are stretched according to the covered expression, as we will see in the next section.

Reversed glyphs, with respect to the vertical — and sometimes also to the horizontal — axis, as in Figure 1, are taken from the Computer Modern font family. For example, there are:

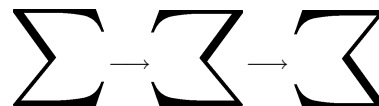
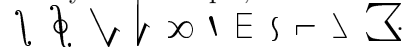


Figure 1: Sum symbol with vertical then horizontal mirror image

Other symbols with mirror image forms already in use¹ are not added to this font. Of course, Latin and Greek alphabetic symbols can be used in Arabic mathematical expressions. In this first phase of the project, we aren't integrating these symbols into the font. They can be brought in from other existing fonts.

3 A Dynamic Font

The composition of variable-sized letters and curvilinear symbols is one of the hardest problems in digital typography. In high-quality printed Arabic works, justification of the line is performed through using the kashida, a curvilinear variable lengthening of letters along the baseline. The composition of curvilinear extensible mathematical symbols is another aspect of dynamic fonts. Here, the distinction between fixed size symbols and those with variable width, length, or with bidimensional variability, according to the mathematical expression covered by the symbol, is of great importance.

Certain systems [11] solve the problem of vertical or horizontal curvilinear extensibility through the a priori production of the curvilinear glyphs for certain sizes. New compositions are therefore necessary beyond these already available sizes. This option doesn't allow a full consideration of the curvilinearity of letters or composed symbols at large sizes. A better approach to get curvilinear letters or extensible mathematical symbols consists of parameterizing the composition procedure of these symbols. The parameters then give the system the required information about the size or the level of extensibility of the symbol to extend. As an example, we will deal with the particular case of the opening and closing parenthesis as *vertically* extensible curvilinear symbol and with the kashida as a *horizontally* extensible curvilinear symbol. This can be generalized to any other extensible symbol.

The CurExt system was developed to build extensible mathematical symbols in a curvilinear way.

¹ The Bidi Mirrored property of characters used in Unicode.

The previous version of this system was able to produce automatically certain dynamic characters, such as parentheses, using METAFONT. In this adaptation, we propose to use the Adobe POSTSCRIPT Type 3 format [13].

The POSTSCRIPT language defines several types of font, 0, 1, 2, 3, 9, 10, 11, 14, 32, 42. Each one of these types has its own conventions to represent and to organize the font information. The most widely used POSTSCRIPT font format is Type 1. However, a dynamic font needs to be of Type 3 [3].

Although the use of Type 3 loses certain advantages of Type 1, such as the possibility of producing hints for when the output device is of low resolution, and in the case of small glyphs, a purely geometrical treatment can't prevent the heterogeneity of characters. Another lost advantage is the possibility of using Adobe Type Manager (ATM) software. These two disadvantages won't arise in our case, since the symbols are generally without descenders or serifs and the font is intended to be used with a composition system such as T_EX, not directly in Windows.

The POSTSCRIPT language [7] produces a drawing by building a path. Here, a path is a set of segments (`lineto`) and third degree Bézier curves (`curveto`). The path can be open or closed on its origin (`closepath`). A path can contain several control points (`moveto`). Once a path is defined, it can be drawn as a line (`stroke`) or filled with a color (`fill`). From the graphical point of view, a glyph is a procedure defined by the standard operators of POSTSCRIPT.

To parameterize the procedure, the form of the glyph has to be examined to determine the different parts of the procedure. This analysis allows determining exactly what should be parameterized. In the case of an opening or closing parenthesis, all the parts of the drawing depend on the size: the width, the length, the boldness and the end of the parenthesis completely depend on the size. Figure 2 shows the variation of the different parameters of the open parenthesis according to the height. We have chosen a horizontally-edged cap with a boldness equal to half of the boldness of the parenthesis. The same process is applied to the kashida.

Producing a dynamic parenthesis such as that in Figure 3 follows these steps:

- collecting the various needed sizes in a parameter file `par`;
- generating a file `p1` with the local tool `par2p1` starting from the `par` file;
- converting the file `p1` into a metric file `tfm` with the application `pltotf`;
- compiling the document to generate a `dvi` file;

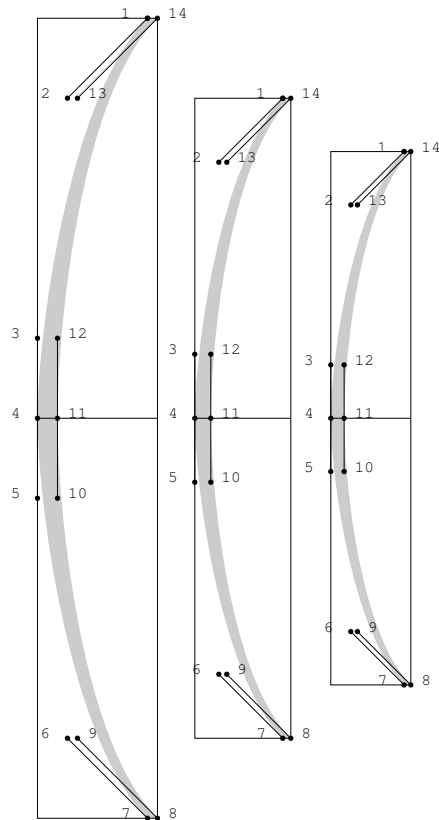


Figure 2: Parametrization of dynamic parenthesis

- converting the file from `dvi` to `ps` format.

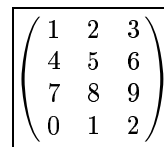
This process should be repeated as many times as needed to resolve overlapping of extensible symbols.

The curvilinear parentheses is produced by CurExt with the following encoding:

```

$ \parentheses{
  \matrix{1 & 2 & 3\cr
          4 & 5 & 6\cr
          7 & 8 & 9\cr
          0 & 1 & 2\cr}
} $

```

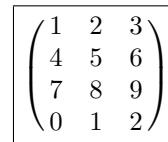


instead of the straight parentheses given by the usual encoding in T_EX:

```

$ \left(
  \matrix{1 & 2 & 3\cr
          4 & 5 & 6\cr
          7 & 8 & 9\cr
          0 & 1 & 2\cr}
\right) $

```

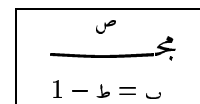


In the same way, we get the curvilinear kashida with CurExt:

```

\amarabmath
${\csum_{b=T-1}^s}$

```



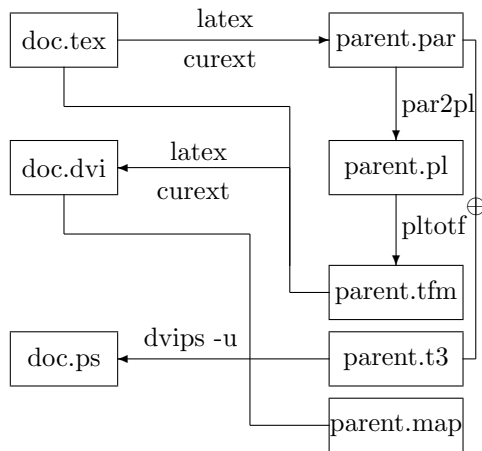
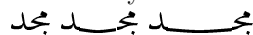


Figure 3: Generation of dynamic parentheses

instead of the straight lengthened one obtained by RyDArab:

```
\amarabmath
 $\sum_{b=T-1}^s$ 
```

We can stretch Arabic letters in a curvilinear way through the kashida by CurExt:



4 Conclusion

The main constraints observed in this work were:

- a close observation of the Arabic calligraphy rules, in the Naskh style, toward their formalization. It will be noticed, though, that we are still far from meeting all the requirements of Arabic calligraphy;
- the heavy use of some digital typography tools, rules and techniques.

RamzArab, the Arabic mathematical font in Naskh style, is currently available as an OpenType font. It meets the requirements of:

- homogeneity: symbols are designed with the same nib. Thus, their shapes, sizes, boldness and other attributes are homogeneous;
- completeness: it contains most of the usual specific Arabic symbols in use.

These symbols are about to be submitted for inclusion in the Unicode standard. This font is under test for Arabic mathematical e-documents [12] after having been structured for Unicode [2, 4].

The dynamic component of the font also works in POSTSCRIPT under CurExt for some symbols such as the open and close parenthesis and the kashida. That will be easily generalized to other variable-sized symbols. The same adaptation can be performed within the OpenType format.

References

- [1] <http://www.linux.org.sa>.
- [2] Jacques André, Caractères numériques : introduction, *Cahiers GUTenberg*, vol. 26, 1997.
- [3] Daniel M. Berry, Stretching Letter and Slanted-baseline Formatting for Arabic, Hebrew and Persian with `ditroff/ffortid` and Dynamic POSTSCRIPT Fonts, *Software-Practice & Experience*, no. 29:15, 1999, pp. 1417–1457.
- [4] Charles Bigelow et Kris Holmes, Création d’une police Unicode, *Cahiers GUTenberg*, vol. 20, 1995.
- [5] Yannis Haralambous, Une police mathématique pour la Société Mathématique de France : le SMF Baskerville, *Cahiers GUTenberg*, vol. 32, 1999, pp. 5–19.
- [6] Yannis Haralambous and John Plaice, Multilingual Typesetting with Ω , a Case Study: Arabic, *Proceedings of the International Symposium on Multilingual Information Processing (Tsukuba)*, 1997, pp. 137–154.
- [7] Adobe Systems Incorporated, *POSTSCRIPT Language Reference Manual*, Second ed., Addison-Wesley, 1992.
- [8] Klaus Lagally, *ArabTEX — Typesetting Arabic with Vowels and Ligatures*, EuroTEX’92 (Prague), 1992.
- [9] Azzeddine Lazrek, Aspects de la problématique de la confection d’une fonte pour les mathématiques arabes, *Cahiers GUTenberg*, vol. 39–40, Le document au XXI^e siècle, 2001, pp. 51–62.
- [10] Azzeddine Lazrek, A package for typesetting arabic mathematical formulas, *Die TEXnische Komödie*, DANTE e.V., vol. 13. (2/2001), 2001, pp. 54–66.
- [11] Azzeddine Lazrek, CurExt, Typesetting variable-sized curved symbols, EuroTEX 2003 preprints, pp. 47–71 (to appear in *TUGboat*).
- [12] Mustapha Eddahibi, Azzeddine Lazrek and Khalid Sami, *Arabic mathematical e-documents*, International Conference on TEX, XML and Digital Typography (TUG 2004, Xanthi, Greece), 2004.
- [13] Włodzimierz Bzyl, The Tao of Fonts, *TUGboat*, vol. 23, 2002, pp. 27–39.
- [14] Thomas W. Phinney, *TrueType, POSTSCRIPT Type 1 & OpenType: What’s the Difference?*, Version 2.00 (2001).

	0	1	2	3	4	5	6	7	8	9
1x										
2x										
3x				!	#	\$	%	&	'	(
4x)	*	+	,	-	.	/	0	1	2
5x	3	4	5	6	7	8	9	:	:	,
6x	=	-	?	@	.	.	:	:	.	.
7x	% ₀₀	%	% _.	٠	١	٢	٣	٤	٥	٦
8x	٧	٨	٩	٠	١				^	'
9x	ا	ب	ج	د	هـ	و	ز	ح	ط	ك
10x	ع	ف	ق	ك	ل	م	ن	هـ	و	ز
11x	ح	ط	ك	ل	م	ن	هـ	و	ز	ح
12x	د	هـ	و	ز	ح	ط	ك	ل	م	ن
13x	س	ع	ف	ص	ق	ح	هـ	ط	ك	ل
14x	ل	م	ن	هـ	و	ز	ح	ط	ك	ل
15x	هـ	و	ز	ح	ط	ك	ل	م	ن	هـ
16x	ع	ف	ص	ق	ك	ل	م	ن	هـ	و
17x	ر	ط	ي	ك	ل	م	ن	س	ع	ف
18x	ص	ق	ك	ل	م	ن	هـ	و	ز	ح
19x	٠	١	٢	٣	٤	٥	٦	٧	٨	٩
20x	١	٢								

Figure 4: RamzArab Arabic mathematical font

Creating Type 1 Fonts from METAFONT Sources: Comparison of Tools, Techniques and Results

Karel Píška

Institute of Physics, Academy of Sciences

182 21 Prague

Czech Republic

piska@fzu.cz

<http://www-hep.fzu.cz/~piska/>

Abstract

This paper summarizes experiences in converting METAFONT fonts to PostScript fonts with `TeXtrace` and `mfttrace`, based on programs of autotracing bitmaps (`AutoTrace` and `potrace`), and with systems using analytic conversion (`MetaFog` and `MetaType1`, using `METAPOST` output or `METAPOST` itself). A development process is demonstrated with public Indic fonts (Devanagari, Malayalam). Examples from the Computer Modern fonts have been also included to illustrate common problems of conversion. Features, advantages and disadvantages of various techniques are discussed. Postprocessing—corrections, optimization and (auto)hinting—or even preprocessing may be necessary, before even a primary contour approximation is achieved. To do fully automatic conversion of a perfect METAFONT glyph definition into perfect Type 1 outline curves is very difficult at best, perhaps impossible.

KEYWORDS: font conversion, bitmap fonts, METAFONT, METAPOST, outline fonts, PostScript, Type 1 fonts, approximation, Bézier curves.

1 Introduction

In recent years, several free programs for creating PostScript outline fonts from METAFONT sources have been developed. The aim of this paper is to give a short comparison of these programs, with references to original sources and documentation, and to provide a brief description of their use. We will discuss advantages and drawbacks, and demonstrate numerous examples to compare important features and to illustrate significant problems. We omit technical details described in the original documentation and concentrate our attention on the quality of the output, including hinting issues.

The programs `TeXtrace` and `mfttrace` read original METAFONT sources, generate high-resolution `pk` bitmaps, call autotracing programs (`AutoTrace` or `potrace`) and finally generate the files in the Type 1 format (`pfb` or `pfa`).

`MetaType1` creates Type 1 output from `METAPOST` sources. Therefore it requires rewriting font definitions from METAFONT into `METAPOST`.

Similarly, `MetaFog` converts the PostScript files generated by `METAPOST` to other PostScript files containing only outlines, that can be subsequently assembled into Type 1 fonts. `MetaFog` is not a new

product, but its excellent results remain, in our comparisons, unsurpassed.

Additionally, we may need adequate encoding files. If none are available, a `TeX` encoding (e.g., the standard `TeX T1` encoding) is usually used as the default.

2 Autotracing Bitmaps

2.1 `TeXtrace` with `AutoTrace`

Péter Szabó developed `TeXtrace` [18]. It is a collection of Unix scripts. It reads the original METAFONT sources, rendering the font bitmaps into PostScript (via `dvips`). For converting the resulting bitmaps to outlines, it calls (in the version of 2001) the `AutoTrace` program [21] created by Martin Weber, and, finally, composites the final files in the Type 1 format. `TeXtrace` works fully automatically and can be invoked by a command like this:

```
bash traceall.sh mfname psname psnumber
where mfname.mf is the name of the METAFONT font, psname.pfb is the name of the Type 1 font file, and psnumber denotes a Type 1 UniqueID [1].
```

The *Adobe Type 1 Font Format* documentation [1, pp. 29–33] recommends observing certain *Type 1*

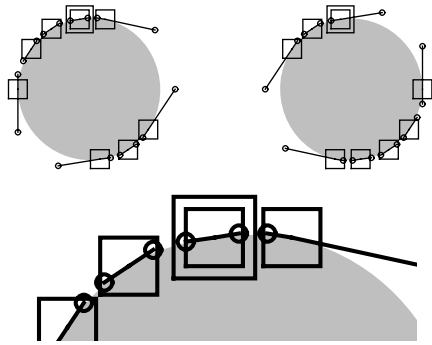


Figure 1: TeXtrace: “¨” in cmr10.

conventions: 1) points at extremes; 2) tangent continuity; 3) conciseness; and 4) consistency.

The outline results from TeXtrace (that is, from AutoTrace) are relatively faithful to the original bitmaps. Some artifacts exist, but they are invisible in usual font sizes and magnifications and for practical purposes may be negligible. Nonetheless, they spoil our attempts to automatically produce perfect, hinted, outline fonts.

The underlying reason is that the information about the control points in the original METAFONT is lost, and the Type 1 conventions are not satisfied, as exemplified in figure 1. The endpoints (double squares) are not placed at extremes (rule 1), most of the horizontal and vertical points of extrema are missing. On the other hand, the outline definition is not concise (rule 3) — due to the large numbers of control points in the glyph definitions, the font files generated by TeXtrace are huge. Furthermore, the two identical periods in the dieresis glyph “¨” are approximated by different point sets (rule 4).

The following examples show the results of conversion of Indic fonts submitted to TUG India 2002 [16], devanagari (dvng10) and Malayalam (mm10). Typical irregularities produced by conversion with TeXtrace are *bumps* and *holes*. Figure 2 demonstrates bumps caused by the envelope being stroked along a path with a rapid change of curvature, and by cases of transition from a straight line to a significantly small arc. The second clipped part of the letter “pha” shows a hole.

I tried to remove those bumps and holes, and (partially) other irregularities at the Type 1 level with a set of special programs manually marking places to be changed in a “raw” text, translated by `t1disasm` and by `t1asm` back after modifications (both programs are from the `t1utils` package [13]), which achieves a better outline approximation, as shown in figure 3. The postprocessing consisted of:

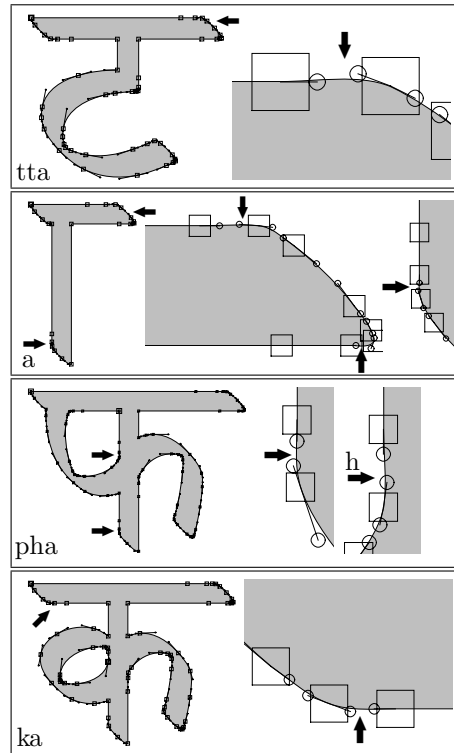


Figure 2: Results of TeXtrace (AutoTrace): bumps and a hole (h).

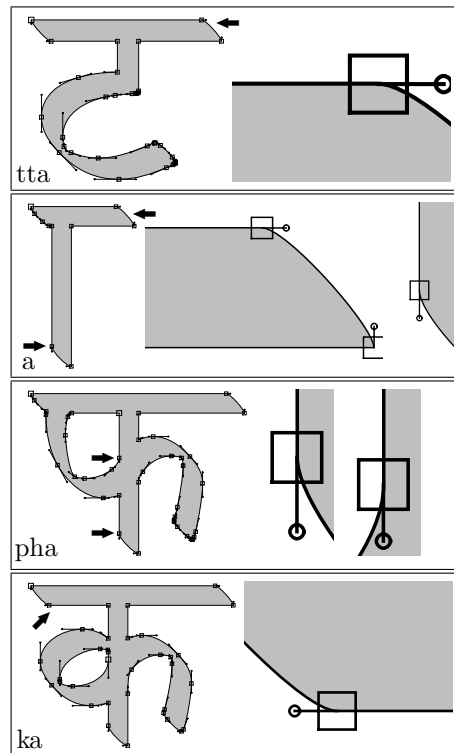


Figure 3: Improved results achieved with postprocessing.

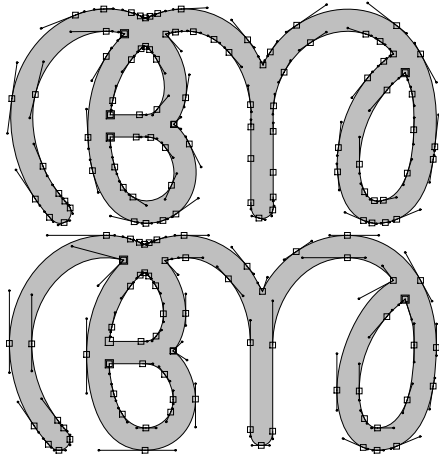


Figure 4: \TeX trace (AutoTrace) first without and then with postprocessing for the Malayalam “a”, showing undetected corners.

inserting missing extrema points, changing the first nodes of contour paths (if desirable), and the optimization of merging pairs (or sequences) of Bézier segments together, and joining nodes in horizontal or vertical straight parts to eliminate redundant nodes.

However, when this process was applied to the Malayalam fonts, we meet another problem: undetected corners in Figure 4. Instead of attempting to correct them, I stopped my postprocessing attempts, and switched to experiments with analytic methods of conversion.

2.1.1 Examples of CM-super

Type 1 fonts [20] generated by Vladimir Volovich (first announced in 2001) inherit typical bugs produced by tracing bitmaps by AutoTrace (as invoked by \TeX trace) such as bumps and holes, improper selection of starting points of contour paths, and problems in distinguishing sharp corners and small arcs. We illustrate them in several following figures, in order to demonstrate that fixing such irregularities automatically is difficult.

In the period “.” from the `sfrm1000` font (its source is the original `cmr10`), an optimization cannot exclude the redundant node (fig. 5) (it is still the starting point of the path).

The minus “-” derived from `cmr10` contains a bump, and minus from `cmtt10` two bumps (fig. 6). Moreover, these bumps have been hinted and have their own hints (probably as results of autohinting).

In the letter “M” from `cmtt10`, we observe missing dishes, a hole and a bad approximation of an arc

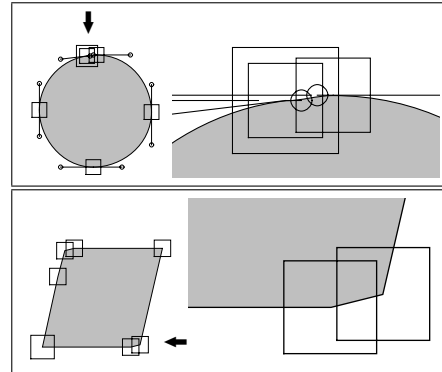


Figure 5: CM-super: period in `sfrm1000` and `sfsi1000`, with redundant node.

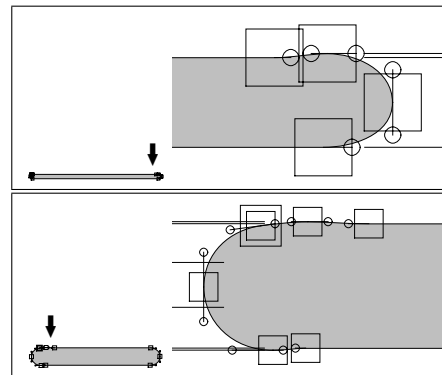


Figure 6: CM-super: minus in `sfrm1000` and `sftt1000`, with bumps.

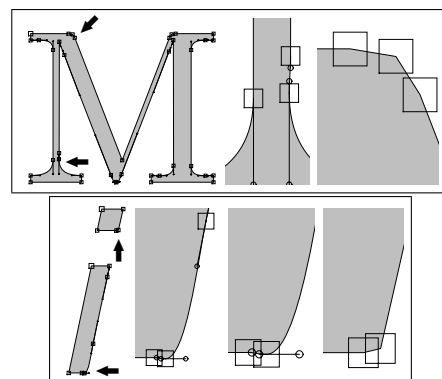


Figure 7: CM-super: “M” in `sfrm1000` and “i” in `sfsi1000`.

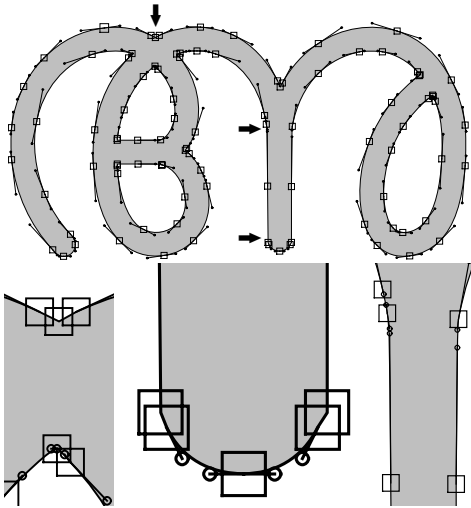


Figure 8: TEX trace (using `potrace`), with different corners.

(fig. 7). On the contrary, in “i” the corners are not detected properly, we also have a hinted bump.

2.2 TEX trace with `potrace`

The 2003 version of TEX trace supports alternative bitmap tracing with `potrace` [17], developed by Peter Selinger. In this version, the real corners are detected or at least detected better than with `AutoTrace` (see fig. 8). Thus, bumps and holes have been suppressed, but smooth connections have often been changed to sharp corners (not present originally). While the bumps demonstrated violation of consistency and may produce invalid hinting zone coordinates (fig. 6), the unwanted sharp corners mean loss of tangent continuity (the middle clip in fig. 8). Unfortunately, the approximation does not preserve horizontal and vertical directions (the right clip), the stem edges are oblique — the difference between the two arrows on the left edge is 2 units in the glyph coordinate space.

2.3 `mftrace`

Han-Wen Nienhuys created `mftrace` [15, 3], a Python script which calls `AutoTrace` or `potrace` (as with TEX trace) to convert glyph bitmap images to outlines. The results of tracing are thus expected to be very similar to those of TEX trace. In fact, for the analyzed Indic fonts, they are identical, as we can see in the first image in figure 9 (compare with TEX trace results in fig. 4). With the `--simplify` option, `mftrace` calls `FontForge` [22] (previously named `PfaEdit`) to execute postprocessing simplification; this helps to exclude redundant nodes from outline contours, as in the second image in figure 9.

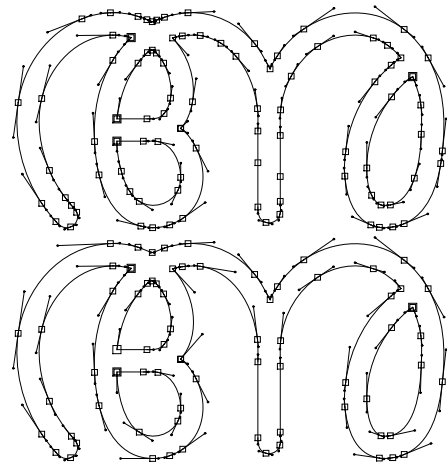


Figure 9: `mftrace` without and with `--simplify`.

3 Analytic Conversions

3.1 `MetaType1`

`MetaType1` [8, 9] is a programmable system for auditing, enhancing and generating Type 1 fonts from METAFONT sources. `MetaType1` was designed by Bogusław Jackowski, Janusz M. Nowacki and Piotr Strzelczyk. The `MetaType1` package is available from <ftp://bop.eps.gda.pl/pub/metatype1> [10].

This “auditing and enhancing” is a process of converting the Type 1 font into `MetaType1` (text) files, generating proof sheets, analysis, making corrections and regenerating modified Type 1 fonts. It is an important tool for checking, verifying and improving existing Type 1 fonts.

`MetaType1` works with the METAFONT language. Therefore the METAFONT font sources must be converted/rewritten into METAFONT. Macro package extensions of METAFONT and other miscellaneous programs provide generation of proper structure of the Type 1 format, evaluate hints (not only the basic outline curves), and create `pfb` and also `afm` and `pfm` files.

During the rewriting process, users define several parameters of the Type 1 font, including the PostScript font encoding — PostScript glyph names and their codes — because METAFONT sources do not contain this data in a form directly usable for Type 1 encoding vectors. METAFONT output commands have to be changed to their METAFONT alternatives. Similarly, it is necessary to substitute METAFONT commands not available in METAFONT, to define METAFONT variants of pen definitions and pen stroking, etc.

Alternative METAFONT commands are defined in the `MetaType1` files `fontbase.mp`, `plain_ex.mp`,

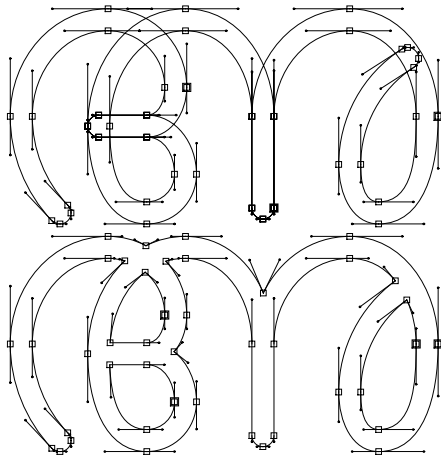


Figure 10: MetaType1 — primary outlines and overlap removal.

et al. Other (new) commands may be defined by the user. Correspondence between METAFONT and METAPOST is approximately as shown in the following table (of course, the details may vary from font to font):

METAFONT	METAPOST
<code>fill path;</code> <code>draw path;</code>	<code>Fill path;</code> <code>pen_stroke() (path) (glyph);</code> <code>Fill glyph;</code>
<code>penlabels(1,2);</code> <code>beginchar(...</code> <code>endchar;</code>	<code>justlabels(1,2);</code> <code>beginglyph(...</code> <code>endglyph;</code>

Many METAFONT commands have no counterpart in METAPOST [6]. For example, operations with bitmap pictures: in METAPOST, font data is represented as PostScript curves, not bitmaps. As a

result, writing METAPOST code that would produce equivalent results as original METAFONT code using these or other such features would be very difficult.

After the basic conversion, the next step is removing overlaps (if any are present) using the MetaType1 command `find_outlines`. Figure 10 shows the results before and after overlap removal for the Malayalam vowel *a* (font `mm10` using pen stroking with a circular pen). This operation is not necessary in METAFONT, since it generates bitmaps. In the METAPOST environment of PostScript outlines, however, we need to reduce overlapping curves to single or pairs of paths.

MetaType1 also allows insertion of commands for automatic computation of horizontal and vertical hints (`FixHStems`, `FixVStems`). The Type 1 font can be visualized in a proof sheet form containing the control point labels (numbers) and hinting zones (figure 11).

So far, so good. But there are two crucial problems. First, the METAFONT Malayalam fonts designed by Jeroen Hellingman [5], use the command

```
currenttransform := currenttransform
                  shifted (.5rm, 0);
```

So all the glyphs should be shifted to the right. METAFONT saves the transformation command and does this operation automatically. By contrast, in METAPOST we need to insert the `shift` commands explicitly in all glyph programs. Also the labels must be shifted! In my experiments, I did this shift operation later, before final assembly of the Type 1 fonts.

The second problem is that in MetaType1 (I used MetaType1 version 0.40 of 2003) a regular pen stroking algorithm is not available, only a simplified method of connecting the points ‘parallel’ to the

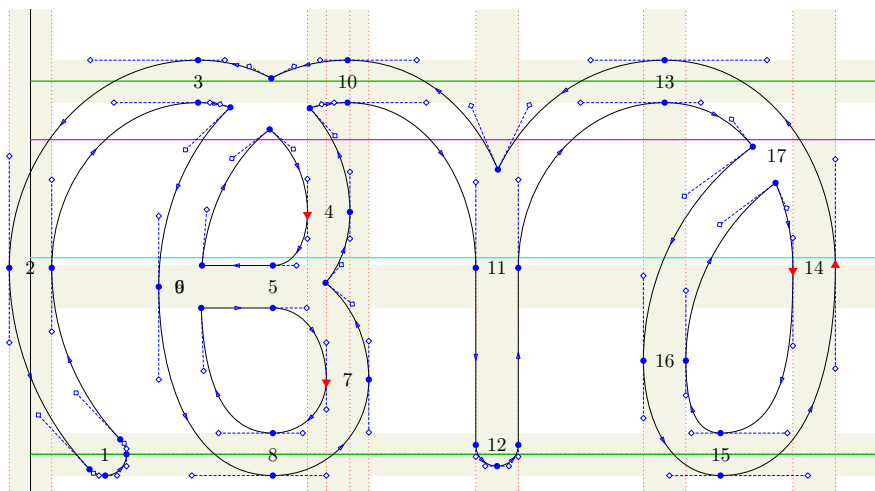


Figure 11. MetaType1 — proof sheet.

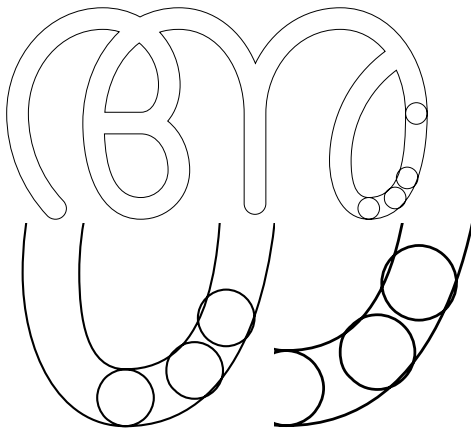


Figure 12: MetaType1 — bad pen stroking.

nodes on the path. Therefore the approximation of the envelope is not correct. For example, in Figure 12 it should be asymmetric, but it is symmetric. Inserting additional nodes cannot help, because the bisection results will again be asymmetric. The figure shows the outline curves do not correspond to the real pen in two midpoint locations. The envelope there looks narrow and it is in fact narrower than it should be. I hope that this problem could be solved in a future release, at least for pen stroking with a circular pen.

Even more serious is a situation with the rotated elliptical pen used in the Devanagari fonts designed by Frans Velthuis [19] (and also other Indic fonts derived from `dvng`). Absence of a regular pen stroking in MetaType1 makes it impractical for such complicated fonts. MetaType1 approximates the pen statically in path nodes, tries to connect their static end points, and ignores complicated dynamic correlations between the path, the pen and the envelope. Unfortunately, in this case the results of the envelope approximation are not correct and cannot be used (figure 13).

3.2 MetaFog

Two programs using analytic conversion were presented in 1995. Basil K. Malyshev created his BaKoMa collection [14] and Richard J. Kinch developed MetaFog [11]. BaKoMa is a PostScript and TrueType version of the Computer Modern fonts. Malyshev’s paper discusses some problems of conversion, especially regarding hinting, but his programs and detailed information about the conversion algorithm are not available.

R. Kinch created MetaFog along with `weeder`, which supports interactive processing of outlines,

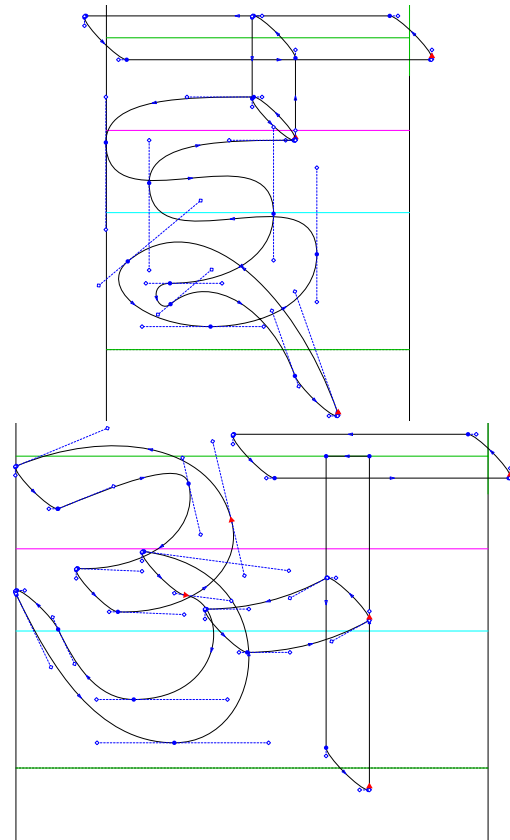


Figure 13: MetaType1 — Devanagari “i”, “a”.

and a package for making final fonts from outlines generated by MetaFog in TrueType, Type 1 and other formats. MetaFog itself (I used an evaluation version graciously donated by Richard) reads the METAPOST output from the command:

```
mpost '&mfplain options;' input fontname.mf
```

Thus, the conversion (from METAFONT sources) is limited to fonts that can be processed by METAPOST, that is, do not contain METAFONT-specific definitions and commands. MetaFog generates another PostScript file consisting only of the outline structures. A conversion process is described also in the paper written by Taco Hoekwater [7].

MetaFog evaluates outline contours and precisely computes envelopes of an elliptical pen stroking along a Bézier curve. We must notice that the envelopes in general are not cubic Bézier curves and their representation in a Type 1 font must be an approximation. The results for a circular pen, on the other hand, can be considered perfect. Figures 14 and 15 show an example of the Malayalam letter “a” (font `mm10`): the initial and final contours and the final Type 1 font with control points (stroked

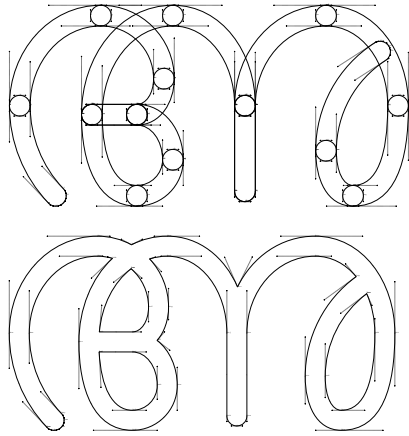


Figure 14: MetaFog—initial input contour and final result.

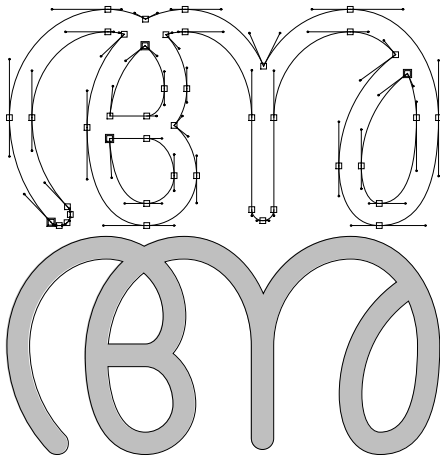


Figure 15: MetaFog—final Type 1 font.

version) and its visual comparison with METAFONT output embedded in a Type 3 font, respectively.

3.2.1 Problems with Complex Pen-stroking

A more complicated situation is the conversion of fonts using pen stroking with a rotated elliptical pen, such as the Devanagari font. Figure 16 illustrates this case. The initial input contour and final result contour (tta1) look good—in the first image we can see the projections of the pen in nodes corresponding to METAFONT source. But exact comparison with the original METAFONT output embedded in a Type 3 font (tta2) and primary MetaFog conversion displayed together with the METAFONT source (tta3) shows that this approximation is not correct. Because these elements are very common in shapes

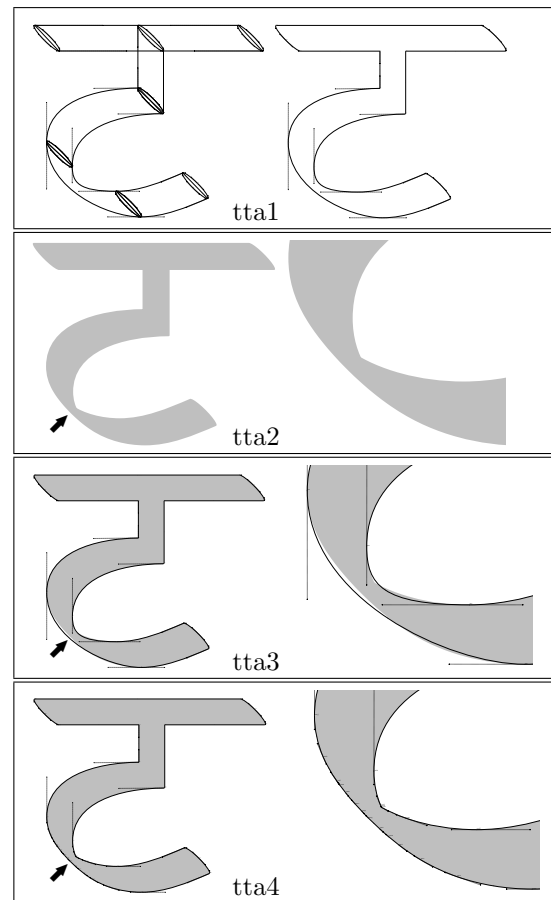


Figure 16: MetaFog contours, METAFONT output, primary and secondary conversion on the METAFONT background.

of all but the simplest Devanagari glyphs, corrections are necessary.

I therefore applied a simple pen-dependent preprocessing step before the MetaFog conversion, thus adapting the METAFONT output as a modified form of bisection, as discussed in a paper by R. Kinch [11]. The preprocessing scans curves, searching for points where the path direction and the direction of main axis of the pen coincide (namely 135°) and inserts these points as additional path nodes. In our case, the transformation matrix is $\cos \theta * [1, 1, -1, 1]$, so we solve only a quadratic equation and can find 0, 1 or 2 (at most) of these points. This technique corrects the MetaFog approximation of all such occurrences in the `dvng` font. The result of this secondary MetaFog conversion with METAFONT source is shown in the last panel of Figure 16 (tta4).

Similar improvements for the Devanagari letters “a” and “pha” are shown in figure 17. For “pha”, the first 135° node was already present

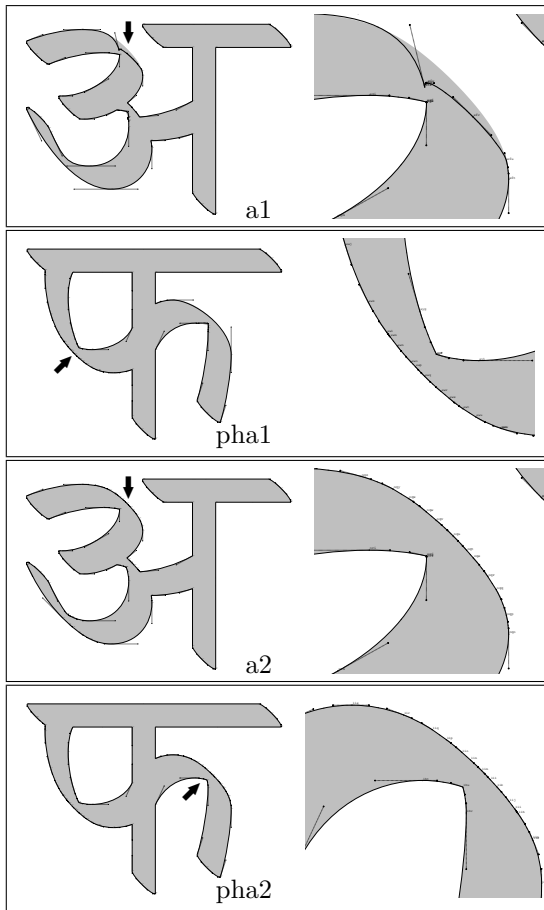


Figure 17: MetaFog output before and after modification of METAPOST source.

in the path defined by the METAFONT source (first panel, pha1); on the contrary, the second occurrence of a 135 degree point was absent, and therefore it was inserted in the METAPOST output (last panel, pha2).

Of course, this improvement is not universal, it only solves a special problem with a special pen for a special font.

Figure 18 illustrates movement of a rotated elliptical pen stepping along a “nice” path (panel 1). However, correlations with the pen are not trivial: changes of curvature of the outer wingtip curve do not have simple monotonic behavior, and the inner wingtip curve (panel 2) is even more complicated. This means that the pen-stroked wingtip curves along a single Bézier curve cannot be approximated by single Bézier curves (compare with the starting fig. 16, panel tta1), i.e., an envelope edge of a pen along a *simple* path is *not simple*.

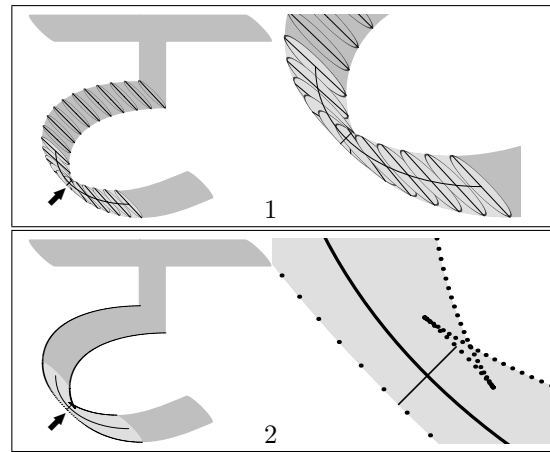


Figure 18: Wingtip curves in METAPOST source.

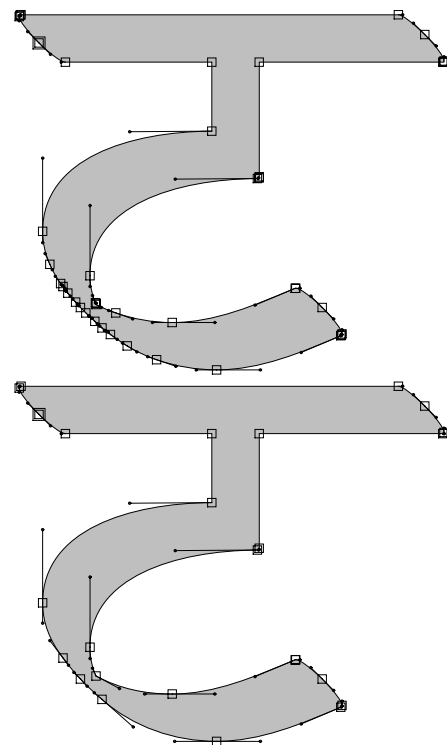


Figure 19: MetaFog converted to Type 1—before and after postprocessing.

3.2.2 Automatic Conversion Problems

A “dark side” of improving the curve approximation is a fragmentation of an envelope curve into many segments (often more than 10, and up to 16 in Devanagari!). We achieve a faithful approximation (limited only by numerical accuracy) at the expense of conciseness. To make up for this, postprocessing is needed. The original MetaFog output and a

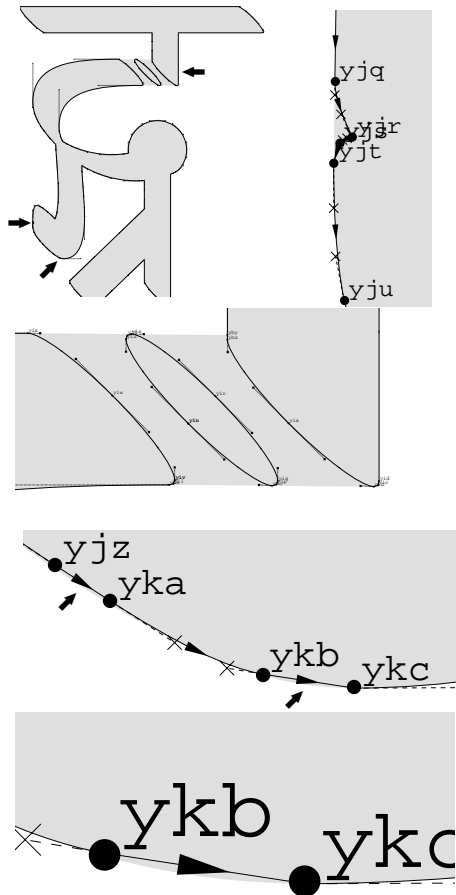


Figure 20: MetaFog — problems with automatic conversion.

result of my (preliminary) optimization assembled into Type 1 fonts is shown in Figure 19.

Unfortunately, even a small computational inaccuracy can make automatic conversion and optimization impossible, and even make it very difficult to design postprocessing algorithms. In Figure 20, we demonstrate problems with the primary approximation of an envelope stroked by a rotated elliptical pen, and also difficulties with automatic optimization of the Devanagari ligature “d+g+r”.

In the first panel of fig. 20, we observe an artifact produced by MetaFog due to a complicated correlation of the pen and the path. Fortunately, those cases are very rare (less than 1% of glyphs in Devanagari).

In the second panel, the path and subsequently the corresponding envelope edges are not absolutely horizontal, thus (probably) MetaFog cannot properly find intersection points and join reconstructed outline components. Those defects are present in

more than 12% of the Devanagari glyphs. In all cases, they have been successfully solved manually by the interactive *weeder* program.

In the last two details in fig. 20 (the lower ending part of the left stem) we can see that both nodes of the left segment are outside the filled area boundary defined by the METAFONT curve. The outer wingtip edge is split there into many segments, some being straight lines — and they should not be, e.g., the first and the third segment marked by 2 arrows in the clip — their curvatures are for us undefined. Additionally, we cannot detect the last segment (magnified in the figure) as horizontal because its angle is “greater than some ϵ ”.

Thus, neither node coordinates, nor segment directions, nor curvatures are reliable. It gives a *visual* comparison of the METAFONT output with its outline approximation. Therefore, my (first and “simple”) idea cannot succeed. This was to classify the behavior of directions and curvatures of *all* the segments *automatically*, and then to divide segments into groups according to directions and curvatures, then *automatically* merging the groups to single Bézier segments. As demonstrated, this optimization may fail or produce incorrect results and, unfortunately, human assistance is needed.

4 Summary

Here we summarize the most important features of the conversion programs found in our experiments.

4.1 Approximate Conversions: *TeXtrace*, *mfttrace*

Advantages:

- approximation covers original METAFONT fonts and correspondence to *pk* bitmaps is (reasonably) close
- simple invocation, robust solution
- fully automatic processing can generate complete, final Type 1 fonts

Disadvantages:

- approximate conversions give only approximate outlines
- lost information about nodes and other control points
- final fonts do not satisfy the *Type 1 conventions*
- *AutoTrace*: problems with recognizing corners, generation of unwanted bumps and holes
- *potrace*: sharp connections, thus loss of tangent continuity, violation of horizontal or vertical directions
- automatic and correct (auto)hinting may yield poor results due to these irregularities

4.2 MetaType1

Advantages:

- complete support for Type 1 font generation
- manual insertion of hinting information possible via simple hinting commands
- font file compression via subroutines

Disadvantages:

- conversion of METAFONT to METAPOST often requires manual rewriting, possibly non-trivial and time-consuming
- bad pen stroking algorithm; in particular, results for complicated fonts using rotated elliptical pens are unusable
- difficulties with removing overlaps in tangential cases

4.3 MetaFog

Advantages:

- fully automatic conversion of METAPOST output to outlines
- “typical” fonts usually achieve perfect results
- even for very complex fonts (again, with rotated elliptical pens), adaptations of METAPOST output and manual editing with *weeder* make it plausible to obtain perfect outlines
- results fulfill the *Type 1 conventions* in most cases (except for those very complex fonts)

Disadvantages:

- MetaFog reads METAPOST output, thus cannot process METAFONT-specific definitions
- complex fonts may still need manual reduction with *weeder* or subsequent optimization of outlines to reach conciseness
- processing is slow

4.4 Final Font Processing and Common Problems

The conversion systems discussed here, with the exception of MetaType1, do not include internal hinting subsystems. To insert hints, we can use font editors, for example FontForge [22]. For successful automatic hinting, however, the font outlines must fulfill certain conditions. Irregularities—absence of nodes at extrema or presence of bumps and holes—are not compatible with autohinting, because extrema points correspond to hinting zones while bumps or holes do not fit them, thus causing outliers. The resulting difference of ± 1 unit in the integer glyph coordinate system, after rounding to integers, is not acceptable for high-quality fonts. Problems may also be caused by other “rounding to

integer” effects, and by the presence of close doublets or triplets.

In my view, these experiments show that the quality of primary outline approximation is crucial to achieve perfect final Type 1 fonts. It is virtually impossible to recreate discarded METAFONT information, or to find exact conditions for a secondary fit that corrects primary contours that were created with irregularities or artifacts. Starting with high-resolution bitmaps is problematic, as too much information has been lost, making subsequent processes of improvement, optimization and hinting difficult at best, not possible to automate and usually not successful.

5 Acknowledgements

I would like to thank all the authors of the free conversion programs, Richard Kinch for donating his MetaFog and *weeder*, the authors of the public METAFONT fonts for Indic languages and other sources used in the contribution, and Karl Berry for help with editing of this article.

References

- [1] Adobe Systems Inc. *Adobe Type 1 Font Format*. Addison-Wesley, 1990.
- [2] Alfred V. Aho, Brian W. Kernighan, and Peter J. Weinberger, *The AWK Programming Language*, Addison-Wesley, 1988.
- [3] Karl Berry. “Making outline fonts from bitmap images.” *TUGboat*, **22**(4), pp. 281–285, 2001.
- [4] Free Software Foundation. GNU awk, <http://www.gnu.org/software/gawk>.
- [5] Jeroen Hellingman. Malayalam fonts, CTAN: language/malayalam.
- [6] John D. Hobby. A User’s Manual for METAPOST. AT&T Bell Laboratories, Computing Science Technical Report 162, 1994.
- [7] Taco Hoekwater. “Generating Type 1 Fonts from METAFONT Sources”, *TUGboat*, **19**(3), pp. 256–266, 1998.
- [8] Bogusław Jackowski, Janusz M. Nowacki, Piotr Strzelczyk. “MetaType1: A METAPOST-based Engine for Generating Type 1 Fonts”, *Proceedings of the XII EuroTEX 2001 conference*, pp. 111–119, Kerkrade, the Netherlands, 23–27 September 2001.
- [9] Bogusław Jackowski, Janusz M. Nowacki, Piotr Strzelczyk. “Programming PostScript Type 1 Fonts Using MetaType1: Auditing, Enhancing, Creating. *Preprints of the XIV EuroTEX 2003*

- conference*, pp. 151–157, Brest, France, 24–27 June 2003 (to appear in *TUGboat*).
- [10] MetaType1distribution: <ftp://bop.eps.gda.pl/pub/metatype1>.
- [11] Richard J. Kinch. “MetaFog: Converting METAFONT Shapes to Contours”, *TUGboat*, **16**(3), pp. 233–243, 1995.
- [12] Donald E. Knuth. *The METAFONTbook*. Addison-Wesley, 1986. Volume C of *Computers and Typesetting*.
- [13] Eddie Kohler. *t1utils* (Type 1 tools), <http://freshmeat.net/projects/t1utils>.
- [14] Basil K. Malyshev, “Problems of the conversion of METAFONT fonts to PostScript Type 1”, *TUGboat*, **16**(1), pp. 60–68, 1995.
- [15] Han-Wen Nienhuys. *mftrace*, <http://www.cs.uu.nl/~hanwen/mftrace>.
- [16] Karel Píška. “A conversion of public Indic fonts from METAFONT into Type 1 format with \TeX -trace.” *TUGboat*, **23**(1), pp. 70–73, 2002.
- [17] Peter Selinger. *potrace*, <http://potrace.sourceforge.net>.
- [18] Péter Szabó. “Conversion of \TeX fonts into Type 1 format”, *Proceedings of the XII Euro-TeX 2001 conference*, pp. 192–206, Kerkrade, the Netherlands, 23–27 September 2001. <http://www.inf.bme.hu/~pts/textrace>; <http://textrace.sourceforge.net>.
- [19] Frans J. Velthuis. Devanagari fonts, CTAN: [language/devanagari](http://ctan.org/language/devanagari).
- [20] Vladimir Volovich. CM-super fonts: CTAN: [fonts/ps-type1/cm-super](http://ctan.org/fonts/ps-type1/cm-super).
- [21] Martin Weber. *AutoTrace*, <http://autotrace.sourceforge.net>.
- [22] George Williams. *FontForge: A PostScript Font Editor*, <http://fontforge.sourceforge.net>.

OpenType and Ω : Past, Present and Future

Yannis Haralambous

Département Informatique

École Nationale Supérieure des Télécommunications de Bretagne

CS 83818, 29238 Brest Cédex, France

yannis.haralambous@enst-bretagne.fr

Gábor Bella

gabor.bella@enst-bretagne.fr

Abstract

This article presents our plans for integrating the OpenType font format into the Ω typesetting system. Beginning with a short summary of what we have achieved so far, we compare potential methods of adaptation. Since most OpenType-related issues we present are valid not only for Ω but for all \TeX -like systems, the authors hope that the article will be interesting for the whole \TeX community.

1 Past

The OpenType font format has been officially available since 1997. In contrast to its predecessors, TrueType and PostScript Type 1 and 2, it provides essential features for proper typesetting of non-LGC¹ scripts, as well as handling LGC ones. Although competing formats with similar capabilities (Apple GX/AAT and Graphite) were and still are available, the marketing force behind OpenType seems strong enough to make it a de facto standard.

The decision to bring OpenType support to Ω as a complement to the already available fonts was taken sometime in 2002 and effective work began around early 2003. At the *EuroTeX 2003* conference, we presented our initial development plans and our first results [1]. Since then, considerable progress has been made, so that in the final version of the EuroTeX proceedings we are finally able to announce a working implementation [1].

However, this implementation is not likely to become *the* final solution, for two reasons. First, it was never intended to be such. Rather, it was created to verify the validity of certain conversion methods and also to provide users with the possibility to install and use OpenType fonts as quickly as possible, even if not all advanced features are available yet. Secondly, as we will see later, complete and robust OpenType support cannot simply be patched onto the existing Ω : as the two do not always share the same philosophy, some parts of Ω will need to

be reorganized in order to take advantage of certain OpenType features.

In the following, we give a quick overview of the present solution and then explain why and how it needs to be revised.

2 Present: OpenType Support, the Quick Way

The word ‘quick’ may seem ironic here: in total, it took us almost a year to produce our first results. However, this was due more to organizational problems than to the difficulty of the task. By ‘the quick way’, we mean the solution that was the most straightforward and the easiest to implement. It is described in detail in the EuroTeX article [1]; here, only a short summary is given.

The present solution is based on the approach that OpenType fonts should be converted to Ω ’s own formats, i.e., OFM (metrics), OVF (virtual fonts) and Ω TP (Ω Translation Process²). Anish Mehta wrote several Python scripts to generate these files, of which the most interesting is perhaps the one that converts the whole OpenType GSUB table (see next section) into Ω TP’s.

Moreover, as the OpenType format is generally not understood by PostScript printers, a conversion to Type 1 (or Type 42) is necessary. To speed up the process, we create Type 1 charstring collections using our own *PFC* tables (see [1] for details) which

¹ Latin-Greek-Cyrillic.

² An Ω TP describes a *finite state automaton* that can be used to filter the input text to perform tasks such as contextual analysis.

are later used by `odvips` to create small, subsetted Type 1 fonts (a.k.a. *minifonts*) on the fly.

The above approach assumes that OpenType can be converted to OFM and Ω TP files without significant loss of information. This is not always the case. In the next section, we will show the main differences between the two formats and why we have decided to abandon some of our results in favour of a better concept. This does not mean that the work done so far was useless: first, we managed to prove that the conversion to Type 1 (which we are planning to keep) is a viable approach and secondly, we are able to provide Ω users with a working, albeit not complete and only temporary, OpenType support.

3 Future: Alternatives of Adaptation

An OpenType font consists mainly of the following parts:

1. font and glyph metrics;
2. Type 2 or TrueType glyph outlines (and hints or instructions);
3. advanced typographic features (mainly GSUB and GPOS);
4. other data tables.

In one form or another, all of them will certainly need to be dealt with either inside Ω or `odvips`. Below, we give a very concise description of our plans regarding each of these fields.

3.1 Metrics

OpenType provides extensive font metric information dispersed among various tables (post, kern, hmtx, hdmx, OS/2, VORG, etc.), both for horizontal and vertical typesetting. In most cases, Ω 's (or \TeX 's) and OpenType's metrics can be converted from one to another, with few but important exceptions (e.g., height/depth, see [1] and [2]). Despite the occasional differences, it seems desirable to dispose of the intermediary OFM files and read OpenType metrics directly from the font file.

3.2 Conversion

As explained in the last section, conversion of OpenType's Type 2 and TrueType outlines to Type 1 has already been implemented. We are also planning to provide Type 42 support for TrueType-flavoured OpenType that would also preserve instructions.

3.3 Advanced features: GSUB and GPOS

These advanced features are the most interesting part of OpenType. The GSUB (glyph substitution) and GPOS (glyph positioning) tables are essential

for typesetting many non-LGC scripts. In Ω , the equivalent of GSUB features are the Ω TP's: they can do everything GSUB features can, including contextual operations. Glyph positioning is a different issue: since the Ω TPs are designed for text rearrangement (substitutions, reordering etc.), they are not suitable for doing glyph placement as easily. In fact, the idea of *microengines*— Ω TP-like Ω plugins—was introduced some time ago, exactly to provide modular, script- and language-specific positioning methods, along the lines of Ω TP files. With the appearance of OpenType fonts, it became clear that positioning features as implemented by the GPOS table and microengines provide essentially the same functionality. It should then be possible to implement microengines as GPOS features and vice versa.

A closely related problem is the fundamental difference between Ω 's and OpenType's way of describing features. Although both use the Unicode encoding, OpenType's GSUB and GPOS features are based on strings of glyph ID's and not of Unicode characters. Ω and some of its Ω TP's, on the other hand, perform tasks such as contextual analysis or hyphenation on character sequences and the passage from characters to 'real' glyph ID's happens only when `(o)dvips` replaces virtual fonts by real ones. Moreover, it is theoretically impossible to convert a glyph-based OpenType feature into a character-based Ω TP, as some glyphs (allographs) may not even have Unicode equivalents.³ The solution to this problem is either to use glyph- and character-based Ω TP's at the same time or else to read GSUB and GPOS features directly from the OpenType font, without conversion to Ω TP's or microengines.

Whichever method we choose, Ω will need to maintain both glyph and character representations of the same text in parallel to be able to perform font-specific (OpenType features) and font-independent (hyphenation) operations at the same time. This dual representation of text is also crucial for the searchability and modifiability of the output (PDF, PS, SVG or any other) document.

3.4 Extensibility

Finally, the OpenType format has the important feature of being extensible: new tables can be added into the font file, containing, for example, data needed by Ω with no OpenType-equivalents (such as metrics or PFC charstrings). Of course, it is necessary that the given font's license allows such additions.

³ There exist some workarounds though, for example to map these glyphs to the Private User Area of Unicode, but this solution is not very elegant, to say the least.

4 Conclusions

It is time to draw conclusions from the arguments made above. First, OpenType seems capable of being a base font format for Ω . Its tabular file structure is flexible enough so that missing Ω -specific information can easily be added to fonts. Metrics can be read directly from the font file (with some exceptions) and then converted on the fly, if necessary. Also, virtual fonts are not useful for OpenType, as it is more reasonable to use a Unicode-based encoding. Thus, neither OFM nor OVF files will be needed when using OpenType fonts. This is one of the arguments against the approach presented in the previous section.

Secondly, in terms of capabilities, Ω TP's and GSUB features are very much compatible. The same can be said for GPOS and microengines, although no implementation exists for the latter yet. However, incompatibility is present on the character/glyph level. As was shown, a dual approach can be helpful, keeping strings of glyph ID's and characters in parallel—another reason why the simple conversion method described in the previous section is not adequate.

In summary, on the one hand we have opted for a tighter integration of Ω and OpenType, by directly reading from (e.g., metrics) and writing to the font file. On the other hand, advanced OpenType features such as glyph positioning and substitution necessitate modifications in Ω 's character handling as well as in the Ω TP/microengine concept. The companion article in this volume will give further details on how we are planning to implement these changes.

References

- [1] Gábor Bella and Anish Mehta: *Adapting Ω to OpenType Fonts*. EuroT \E X 2003 Proceedings. Final version to appear in *TUGboat*.
- [2] Yannis Haralambous and John Plaice: *Ω and OpenType Fonts*. Kyōto University 21st Century COE Program, 2003.
- [3] The OpenType Specification v1.4.
<http://www.microsoft.com/typography/otspec/default.htm>

Moving Ω to a C++-based Platform

John Plaice

School of Computer Science and Engineering
The University of New South Wales
UNSW SYDNEY NSW 2052, Australia
plaice@cse.unsw.edu.au

Paul Swoboda

pswoboda@cse.unsw.edu.au

Abstract

The code for the Omega Typesetting System has been substantially reorganised. All fixed-size arrays implemented in Pascal Web have been replaced with interfaces to extensible C++ classes. The code for interaction with fonts and Omega Translation Processes (OTPs) has been completely rewritten and placed in C++ libraries, whose methods are called by the typesetting engine. The Pascal Web part of Omega no longer uses change files. The overall Omega architecture is now much cleaner than that of previous versions.

1 Introduction

Since the first paper on Omega was presented at the 1993 Aston TUG Conference, numerous experiments have been undertaken in the realm of multilingual typesetting and document processing. This overall work has given important insights into what a future document processing system, including high quality typesetting, should look like. See, for example, the articles published in the 2003 EuroTEX and TUG conferences. Clearly, building an extensive new system will require substantial effort and time, both at the design and the implementation levels. In this article, we present the interim solution for the stabilisation of the existing Omega code base, with a view towards preparing for the design and implementation of a new system.

The standard `web2c` infrastructure, which assumes that a binary is created from a single Pascal Web file and a single Pascal Web change file, is simply not well suited for the development of large scale software, of any genre. For this reason, we have eliminated the change files, and broken up the Pascal Web file into chapter-sized files. All fixed-size arrays have been reimplemented in C++ using the Standard Template Library. Characters are now 32 bits, using the `wchar_t` data type, and character set conversion is done automatically using the routines available in the `stdc++` library. The entire Pascal Web code for fonts and OTPs, including that of Donald Knuth, has been completely rewritten in C++ and placed in libraries. Clean interfaces have

been devised for the use of C++ code from the remaining Pascal code.

This exercise serves two purposes. The first is to stabilise and to simplify the existing Omega distribution. The second is to lay the groundwork for a forthcoming, complete reimplementing of open typesetting software.

2 Problems with Pascal Web

When we examine the difficulties in creating Omega as a derivation of `tex.web`, we should understand that there is no single source for these difficulties.

Pascal was designed so that a single-pass compiler could transform a monolithic program into a running executable. Therefore, all data types must be declared before global variables; in turn, all variables must be declared before subroutines, and the main body of code must follow all declarations. This choice sacrificed ease of programming for ease of compiler development; the resulting constraints can be felt by anyone who has attempted to seriously modify the TEX engine.

Pascal Web attempts to alleviate this draconian language vision by allowing the arbitrary use within code blocks — called *modules* — of pointers to other modules, with a call-by-name semantics. The result is a programming environment in which the arbitrary use of GOTOs throughout the code is encouraged, more than ten years after Dijkstra's famous paper. Knuth had responded correctly to Dijkstra's paper, stating that the reasonable use of

GOTOs simplifies code. However, the arbitrary use of GOTOs across a program, implicit in the Pascal Web methodology, restricts code scalability. Knuth himself has stated that one of the reasons he stopped working on \TeX was his fear that he might break it.

For a skilled, attentive programmer, developing a piece of code that is not going to evolve, it is possible to write working code, *up to a certain level of complexity*. However, for a program that is to evolve significantly, this approach is simply not tenable, because the monolithic Pascal vision is inherited in the change file mechanism of Pascal Web. Modifications to \TeX are supposed to be undertaken solely using change files; the problem with this approach is that the vision of the code maintainer is that they are modifying functions, procedures, and so on. However, the *real structure* of a Pascal Web program is the interaction between the Pascal Web *modules*, not the functions and procedures that they define. As a result, maintaining a Pascal Web program is a very slow process. Back in 1993, when the first Omega work was being undertaken, “slow” did not just mean slow in design and programming, but also in compilation. The slightest modification required a 48-minute recompilation.

The size limitations created by compile-time fixed-size arrays are obvious and well known. This was addressed publicly by Ken Thompson in the early 1980s and has been addressed in Omega as well as in the `web2c` framework, simply by changing these sizes. However, there are other problems in the way these arrays are used. The `eqtb`, `str_pool`, `font_info` and `mem` arrays all have programming interfaces using fixed-size arrays. Whenever these interfaces are insufficient, the \TeX code simply makes direct accesses to these arrays. As a result, any attempt to significantly modify these basic data structures requires the modification of the *entire* \TeX engine, and not simply the implementations of the structural interfaces.

The single input buffer for all active files of `tex.web` turned out to be truly problematic for implementing Omega’s OTPs. Since an OTP can read in an arbitrary amount of text before processing it, a new input buffer had to be introduced to do this collection. The resulting code is anything but elegant, and could certainly be made more efficient.

Finally, problems arise from the `web2c` implementation of Pascal Web. Many of the routines written in C to support the `web2c` infrastructure make the implicit assumption that all characters are 8 bits, making it difficult to generalise to Unicode, even though C itself has a datatype called `wchar_t`.

3 Suitability of C++

The advantages of the use of C++ as an implementation language for stream-oriented typesetting, over the Pascal Web architecture, are manifold. The chief reason for this is that the rich set of tools and methodologies that have evolved in the twenty-five years since the introduction of \TeX includes developments not only in programming languages and environments, but in operating systems, file structure, multiprocessing, and in the introduction of whole new paradigms, including object-oriented software and generic programming.

The STL offers built-in support for arbitrary generic data structures and algorithms, including extensible, random-access arrays. It would be foolish to ignore such power when it is so readily available.

The Standard C++ library also provides built-in wide character support, including, in the GNU `stdc++` implementation, the full `iconv` functionality for character-set conversion between Unicode and any other imaginably-used character set.

C++ is the *de facto* standard for object-oriented systems development, with its capability to provide low-level C-style access to data structures and system resources (and, in the case of Unix-like systems, direct access to the kernel system call API), for the sake of efficiency.

Since C++ is fully compatible with C, one can still take advantage of many existing libraries associated with \TeX , such as Karl Berry’s `kpathsea` file searching library.

The abilities to use well-known design patterns for generic algorithm support (plug-in paragraphers, generic stream manipulation), as well as generic representation of typesetting data itself, add a wealth of possibilities to future, open typesetting implementations.

4 Organisation of the Omega 2 Code Base

Obviously, we are moving on. Our objective is to include the existing Omega functionality, to stretch it where appropriate, leaving clean interfaces so that, if others wish to modify the code base, they can do so. Our current objective is not to rewrite \TeX .

4.1 Pascal Web Components

The `tex.web` file has been split into 55 files called `01.web` to `55.web`. The `tex.ch` file has been converted into 55 files, `01.ch` to `55.ch`. Data structure by data structure, we have passed through the code, throwing out the definitions of the data structures and replacing their uses with Pascal procedure calls

which, once passed through the `web2c` processor, become C++ method calls. In the process, most of the code in the change files ends up either being thrown out, or directly integrated in the corresponding `.web` files.

4.2 Characters, Strings and Files

For characters, \TeX has two types, `ASCII_code` and `text_char`, the respective internal and external representations of 8-bit characters. The new Omega uses the standard C/C++ data type, `wchar_t`. On most implementations, including GNU C++, `wchar_t` is a 32-bit signed integer, where the values `0x0` to `0x7fffffff` are used to encode characters, and the value `0xffffffff` (-1) is used to encode EOF. Pascal Web strings, as interpreted by the `tangle` program, are each assigned a `str_number`, where values 0 to 255 are reserved for the 256 8-bit characters. We have modified `tangle` so that the strings are numbered -256 downwards, rather than 256 upwards. Hence, `str_number` and `wchar_t` can be of the same data type.

When dealing with files, there are two separate issues, the file names, and the file content. Internally, all characters are 4-byte integers, but on most systems, file names are stored using 8-bit encodings, specified according to the user's locale. Hence, character-set conversion is now built into the file-opening mechanisms, be they for reading or writing.

The actual content of the files may come from anywhere in the world and a single file system may include files encoded with many different encoding schemes. We provide the means for opening a file with a specified encoding, as well as opening a file with automatic character encoding detection, using a one-line header at the beginning of the file. The actual character set conversion is done using the `stdc++` local routines. As a result of these choices, the vast majority of the Omega code can simply assume that characters are 4-byte Unicode characters.

4.3 Fonts and OTPs

In terms of numbers of lines written, this is the most significant part of the new Omega; however, because we are using standard OO technology, it is also the most straightforward.

With respect to fonts, significant energy expended, both in the original code as well as in previous Omega implementations, for bit packing of fields in binary font formats, which are stored in memory as they are on disk. By providing a simple OO interface in the character-level typesetter of the Omega engine, we have been able to greatly simplify the code for both the typesetter, as well as the font utilities for conversion between formats.

Similarly, for the OTPs, filters can be implemented as function objects over streams using iterators, tremendously simplifying the code base.

5 Prospects

At the time we are writing, this work is not completely finished. Nevertheless, it is well advanced and detailed documentation will be forthcoming on the Omega website.

If we view things in the longer term, we are clearly moving forward with two related goals, the stabilisation of existing Omega infrastructure, and abandonment of the \TeX infrastructure for the design and implementation of a next-generation open typesetting suite.

Such a suite should be a generic framework with an efficient C++ core, that is universally extensible through a number of well-known scripting interfaces, for example, Perl, Python, and Guile. Implementation of libraries similar to the popular \LaTeX suite could then be done directly in C++, on top of the core API, or as a linked-in C++ stream filter.

Typesetting CJK Languages with Ω

Jin-Hwan Cho
Korean \TeX Users Group
chofchof@ktug.or.kr

Haruhiko Okumura
Mie University
Faculty of Education
514-8507
Japan
okumura@acm.org

Abstract

This paper describes how to typeset Chinese, Japanese, and Korean (CJK) languages with Omega, a 16-bit extension of Donald Knuth's \TeX . In principle, Omega has no difficulty in typesetting those East Asian languages because of its internal representation using 16-bit Unicode. However, it has not been widely used in practice because of the difficulties in adapting it to CJK typesetting rules and fonts, which we will discuss in the paper.

1 Introduction

Chinese, Japanese, and Korean (CJK) languages are characterized by multibyte characters covering more than 60% of Unicode. The huge number of characters prevented the original 8-bit \TeX from working smoothly with CJK languages. There have been three methods for supporting CJK languages in the \TeX world up to now.

The first method, called the *subfont scheme*, splits CJK characters into sets of 256 characters or fewer, the number of characters that a \TeX font metric file can accommodate. Its main advantage lies in using 8-bit \TeX systems directly. However, one document may contain dozens of subfonts for each CJK font, and it is quite hard to insert glue and kerns between characters of different subfonts, even those from the same CJK font. Moreover, without the help of a DVI driver (e.g., *DVIPDFMx* [2]) supporting the subfont scheme, it is not possible to generate PDF documents containing CJK characters that can be extracted or searched. Many packages are based on this method; for instance, *CJK-L^AT_EX*¹ by Werner Lemberg, *HL^AT_EX*² by Koaunghi Un, and the Chinese module in *ConT_EXt*³ by Hans Hagen.

On the other hand, in Japan, the most widely used \TeX -based system is *p \TeX* [1] (formerly known as ASCII Nihongo \TeX), a 16-bit extension of \TeX

localized to the Japanese language. It is designed for high-quality Japanese book publishing (the “p” of *p \TeX* stands for publishing; the name *j \TeX* was used by another system). *p \TeX* can handle multibyte characters natively (i.e., without resorting to subfonts), and it can typeset both horizontally and vertically within a document. It is upward compatible⁴ with \TeX , so it can be used to typeset both Japanese and Latin languages, but it cannot handle Chinese and Korean languages straightforwardly. *p \TeX* supports three widely-used Japanese encodings, JIS (ISO-2022-JP), Shift JIS, and EUC-JP, but not Unicode-based encodings such as UTF-8.

The third route, Omega [3], is also a 16-bit extension of \TeX , having 16-bit Unicode as its internal representation. In principle, Omega is free from the limitations mentioned above, but thus far there is no thorough treatment of how it can be used for professional CJK typesetting and how to adapt it to popular CJK font formats such as TrueType and OpenType. We set out to fill in this blank.

2 CJK Typesetting Characteristics

Each European language has its own hyphenation rules, but their typesetting characteristics are overall fairly similar. CJK languages differ from European languages in that there are no hyphenation

¹ Available on CTAN as `language/chinese/CJK/`

² Available on CTAN as `language/korean/HLaTeX/`

³ Available on CTAN as `macros/context/`

⁴ Although *p \TeX* doesn't actually pass the `trip` test, it is thought to be upward compatible with \TeX in virtually all practical situations.

rules. All CJK languages allow line breaking almost anywhere, without a hyphen. This characteristic is usually implemented by inserting appropriate glues between CJK characters.

One fine point is the treatment of blank spaces and end-of-line (EOL) characters. Korean uses blank spaces to separate words, but Chinese and Japanese rarely use blank spaces. An EOL character is converted in \TeX to a blank space and then to a skip, which is unnecessary for Chinese and Japanese typesetting. To overcome this problem, $\text{p}\TeX$ ignores an EOL when it follows a CJK character.

Moreover, whereas Korean uses Latin punctuation marks (periods, commas, etc.), Chinese and Japanese use their own punctuation symbols. These CJK punctuation symbols need to be treated somewhat differently from ordinary characters. The appropriate rules are described in this paper.

3 CJK Omega Translation Process

We introduce here the CJK Omega Translation Process (CJK- Ω TP)⁵ developed by the authors to implement the CJK typesetting characteristics mentioned above.

An Omega Translation Process (Ω TP) is a powerful preprocessor, which allows text to be passed through any number of finite state automata, which can achieve many different effects. Usually it is quite hard or impossible to do the same work with other \TeX -based systems.

For each CJK language, the CJK- Ω TP is divided into two parts. The first Ω TP (`boundCJK.otp`) is common to all CJK languages, and controls the boundaries of blocks consisting of CJK characters and blank spaces. The second Ω TP (one of `interCHN.otp`, `interJPN.otp`, and `interKOR.otp`) is specific to each language, and controls typesetting rules for consecutive CJK characters.

4 Common Typesetting Characteristics

The first task of `boundCJK.otp` is to split the input stream into CJK blocks and non-CJK blocks, and insert glue (`\boundCJKglue`) in between to allow line breaking.

However, combinations involving some Latin and CJK symbols (quotation marks, commas, periods, etc.), do not allow line breaking. In this case, `\boundCJKglue` is not inserted so that the original line breaking rule is applied. This corresponds to $\text{p}\TeX$'s primitives `\xspace` and `\inhibitxspace`.

`boundCJK.otp` defines seven character sets; the role of each set is as follows.

1. `{CJK}` is the set of all CJK characters; its complement is denoted by `^{\CJK}`.
2. `{XSPACE1}` (e.g., `{‘}`) is the subset of `^{\CJK}` such that `\boundCJKglue` is inserted only between `{CJK}` and `{XSPACE1}` in this order.
3. `{XSPACE2}` (e.g., `)}’;`) is the subset of `^{\CJK}` such that `\boundCJKglue` is inserted only between `{XSPACE2}` and `{CJK}` in this order.
4. `{XSPACE3}` (e.g., `0-9 A-Z a-z`) is the subset of `^{\CJK}` such that `\boundCJKglue` is inserted between `{CJK}` and `{XSPACE3}`, irrespective of the order.
5. `{INHIBITXSPACE0}` (e.g., `—…¥`) is the subset of `{CJK}` *not* allowing `\boundCJKglue` between `{INHIBITXSPACE0}` and `^{\CJK}`, irrespective of the order.
6. `{INHIBITXSPACE1}` (e.g., `、。〉》」』】`), CJK right parentheses and periods) is the subset of `{CJK}` *not* allowing `\boundCJKglue` between `^{\CJK}` and `{INHIBITXSPACE1}` in this order.
7. `{INHIBITXSPACE2}` (e.g., `〈〈「『【【【【`), CJK left parentheses) is the subset of `{CJK}` *not* allowing `\boundCJKglue` in between `{INHIBITXSPACE2}` and `^{\CJK}` in this order.

The second task of `boundCJK.otp` is to enclose each CJK block in a group `{\selectCJKfont_...}`, and convert all blank spaces inside the block to the command `\CJKspace`.

The command `\selectCJKfont` switches to the appropriate CJK font, and `\CJKspace` is defined to be either a `\space` (for Korean) or `\relax` (for Chinese and Japanese) according to the selected language.

Note that if the input stream starts with blank spaces followed by a CJK block or ends with a CJK block followed by blank spaces, then these spaces must be preserved regardless of the language, because of math mode:

```
{CJK} {SPACE} $...$ {SPACE} CJK}}
```

and restricted horizontal mode:

```
\hbox{{SPACE} {CJK} {SPACE}}
```

5 Language-dependent Characteristics

The line breaking mechanism is common to all of the language-dependent Ω TPs (`interCHN.otp`, `interJPN.otp`, and `interKOR.otp`). The glue `\interCJKglue` is inserted between consecutive CJK characters, and its role is similar to the glue `\boundCJKglue` at the boundary of a CJK block.

⁵ Available at <http://project.ktug.or.kr/omega-cjk/>

Some combinations of CJK characters do not allow line breaking. This is implemented by simply inserting a `\penalty 10000` before the relevant `\interCJKglue`. In the case of `boundCJK.otp`, however, no `\boundCJKglue` is inserted where line breaking is inhibited.

The CJK characters not allowing line breaking are defined by the following two classes in `interKOR.otp` for Korean typesetting.

1. `{CJK_FORBIDDEN_AFTER}` does not allow line breaking between `{CJK_FORBIDDEN_AFTER}` and `{CJK}` in this order.
2. `{CJK_FORBIDDEN_BEFORE}` does not allow line breaking in between `{CJK}` and `{CJK_FORBIDDEN_BEFORE}` in this order.

On the other hand, `interJPN.otp` defines six classes for Japanese typesetting, as discussed in the next section.

6 Japanese Typesetting Characteristics

Most Japanese characters are designed on a square ‘canvas’. $\text{p}\text{T}\text{E}\text{X}$ introduced a new length unit, `zw` (for *zenkaku* width, or full-width), denoting the width of this canvas. The CJK- Ω TP defines `\zw` to denote the same quantity.

For horizontal (left-to-right) typesetting mode, the baseline of a Japanese character typically divides the square canvas by 0.88 : 0.12. If Japanese and Latin fonts are typeset with the same size, Japanese fonts appear larger. In the sample shown in Figure 1, Japanese characters are typeset 92.469 percent the size of Latin characters, so that 10 pt (1 in = 72.27 pt) Latin characters are mixed with 3.25 mm (= 13 Q; 4 Q = 1 mm) Japanese characters. Also, Japanese and Latin words are traditionally separated by about 0.25 `zw`, though this space is getting smaller nowadays.

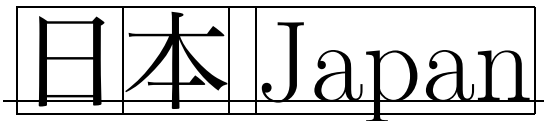


Figure 1: The width of an ordinary Japanese character, 1 `zw`, is set to 92.469% the design size of the Latin font, and a gap of 0.25 `zw` is inserted. The baseline is set to 0.12 `zw` above the bottom of the enclosing squares.

Some characters (such as punctuation marks and parentheses) are designed on a half-width canvas: its width is 0.5 `zw`. For ease of implementation, actual glyphs may be designed on square canvases.

We can use the virtual font mechanism to map the logical shape and the actual implementation.

`interJPN.otp` divides Japanese characters into six classes:

1. Left parentheses: ‘ “ ({ [{ { < < 「 『 【
Half width, may be designed on square canvases flush right. In that case we ignore the left half and pretend they are half-width, e.g., `\hbox to 0.5zw{\hss}`. If a class-1 character is followed by a class-3 character, then an `\hskip 0.25zw minus 0.25zw` is inserted in between.
2. Right parentheses: \ , ’ ”) }] } > > 』 』
Half width, may be designed flush left on square canvases. If a class-2 character is followed by a class-0, -1, or -5 character, then an `\hskip 0.5zw minus 0.5zw` is inserted in between. If a class-2 character is followed by a class-3 character, then a `\hskip 0.25zw minus 0.25zw` is inserted in between.
3. Centered points: ∙ ∙ ∙ ;
Half width, may be designed centered on square canvases. If a class-3 character is followed by a class-0, -1, -2, -4, or -5 character, then an `\hskip 0.25zw minus 0.25zw` is inserted in between. If a class-3 character is followed by a class-3 character, then an `\hskip 0.5zw minus 0.25zw` is inserted in between.
4. Periods: 。 。 。
Half width, may be designed flush left on square canvases. If a class-4 character is followed by a class-0, -1, or -5 character, then an `\hskip 0.5zw` is inserted in between. If a class-4 character is followed by a class-3 character, then an `\hskip 0.75zw minus 0.25zw` is inserted in between.
5. Leaders: —……
Full width. If a class-5 character is followed by a class-1 character, then an `\hskip 0.5zw minus 0.5zw` is inserted in between. If a class-5 character is followed by a class-3 character, then an `\hskip 0.25zw minus 0.25zw` is inserted in between. If a class-5 character is followed by a class-5 character, then a `\kern 0zw` is inserted in between.
0. Class-0: everything else.
Full width. If a class-0 character is followed by a class-1 character, then an `\hskip 0.5zw minus 0.5zw` is inserted in between. If a class-0 character is followed by a class-3 character, then an `\hskip 0.25zw minus 0.25zw` is inserted in between.

Chinese texts can be typeset mostly with the same rules. An exception is the comma and the

period of Traditional Chinese. These two letters are designed at the center of the square canvas, so they should be treated as Class-3 characters.

7 Example: Japanese and Korean

Let us discuss how to use CJK-ΩTP in a practical situation. Figure 2 shows sample output containing both Japanese and Korean characters, which is typeset by Omega with the CJK-ΩTP and then processed by DVIPDFMx.

TeX はスタンフォード大学のクヌース教授によって開発された組版システムであり、組版の美しさと強力なマクロ機能の特徴としている。

TeX은 스탠포드 대학의 크누스 교수에 의해 개발된組版 시스템으로,組版의美와強力한 매크로 기능이 특징이다.

Figure 2: Sample CJK-ΩTP output.

The source of the sample above was prepared with the text editor Vim as shown in Figure 3. Here, the UTF-8 encoding was used to see Japanese and Korean characters at the same time. Note that the backslash character (\) is replaced with the yen currency symbol in Japanese fonts.

```

¥input omega-cjk-sample
¥hsize=75mm ¥par indent=¥zw
{¥japanese
¥TeXはスタンフォード大学のクヌース教授によって
開発された組版システムであり、組版の美しさと強
力なマクロ機能の特徴としている。
}
¥par¥vskip 10pt
{¥korean
¥TeX은 스탠포드 대학의 크누스 교수에 의해 개발
된組版 시스템으로,組版의美와強力한 매크로
기능이 특징이다.
}
¥bye

```

Figure 3: Sample CJK-ΩTP source.

The first line in Figure 3 calls another TeX file omega-cjk-sample.tex which starts with the following code, which loads⁶ the CJK-ΩTP.

```

\ocp\OCPindefault=inutf8
\ocp\OCPboundCJK=boundCJK
\ocp\OCPinterJPN=interJPN
\ocp\OCPinterKOR=interKOR

```

⁶ Omega requires the binary form of ΩTP files compiled by the utility otp2ocp included in the Omega distribution.

Note that inutf8.otp has to be loaded first to convert the input stream encoded with UTF-8 to UCS-2, the 16-bit Unicode.

```

\ocplist\CJKOCP=
  \addafterocplist 1 \OCPboundCJK
  \addafterocplist 1 \OCPindefault
  \nullocplist
\ocplist\JapaneseOCP=
  \addbeforeocplist 2 \OCPinterJPN \CJKOCP
\ocplist\KoreanOCP=
  \addbeforeocplist 2 \OCPinterKOR \CJKOCP

```

The glues \boundCJKglue and \interCJKglue for CJK line breaking mechanism are defined by new skip registers to be changed later according to the language selected.

```

\newskip\boundCJKskip % defined later
\def\boundCJKglue{\hskip\boundCJKskip}
\newskip\interCJKskip % defined later
\def\interCJKglue{\hskip\interCJKskip}

```

Japanese typesetting requires more definitions to support the six classes defined in interJPN.otp.

```

\newdimen\zw \zw=0.92469em
\def\halfCJKmidbox#1{\leavevmode%
  \hbox to .5\zw{\hss #1\hss}}
\def\halfCJKleftbox#1{\leavevmode%
  \hbox to .5\zw{#1\hss}}
\def\halfCJKrightbox#1{\leavevmode%
  \hbox to .5\zw{\hss #1}}

```

Finally, we need the commands \japanese and \korean to select the given language. These commands have to include actual manipulation of fonts, glues, and spaces.

```

\font\defaultJPNfont=omrml
\def\japanese{%
  \clearocplists\pushocplist\JapaneseOCP
  \let\selectCJKfont\defaultJPNfont
  \let\CJKspace\relax % remove spaces
  \boundCJKskip=.25em plus .15em minus .06em
  \interCJKskip=0em plus .1em minus .01em
}
\font\defaultKORfont=omhysm
\def\korean{%
  \clearocplists\pushocplist\KoreanOCP
  \let\selectCJKfont\defaultKORfont
  \let\CJKspace\space % preserve spaces
  \boundCJKskip=0em plus .02em minus .01em
  \interCJKskip=0em plus .02em minus .01em
}

```

It is straightforward to extend these macros to create a L^AT_EX (Lambda) class file.

8 CJK Font Manipulation

At first glance, the best font for Omega seems to be the one containing all characters defined in 16-bit Unicode. In fact, such a font cannot be constructed.

There are several varieties of Chinese letters: Traditional letters are used in Taiwan and Korea, while simplified letters are now used in mainland China. Japan has its own somewhat simplified set. The glyphs are significantly different from country to country.

Unicode unifies these four varieties of Chinese letters into one, if they look similar. They are *not* identical, however. For example, the letter ‘bone’ has the Unicode point 9AA8, but the top part of the Chinese Simplified letter and the Japanese letter are almost mirror images of each other, as shown in Figure 4. Less significant differences are also distracting to native Asian readers. The only way to overcome this problem is to use different CJK fonts according to the language selected.

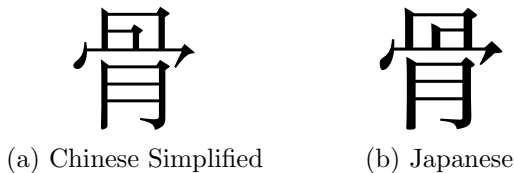


Figure 4: Two letters with the same Unicode point.

OpenType (including TrueType) is the most popular font format for CJK fonts. However, it is neither easy nor simple, even for \TeX experts, to generate OFM and OVF files from OpenType fonts.

The situation looks simple for Japanese and Chinese fonts, having fixed width, because one (virtual) OFM is sufficient which can be constructed by hand. However, Korean fonts have proportional width. Since most of the popular Korean fonts are in OpenType format, a utility that extracts font metrics from OpenType fonts is required.

There are two patches of the `ttf2tfm` and `ttf2pk` utilities⁷ using the `freetype` library. The first,⁸ written by one of the authors, Jin-Hwan Cho, generates OFM and OVF files from TrueType fonts (not OpenType fonts). The other,⁹ written by Won-Kyu Park, lets `ttf2tfm` and `ttf2pk` run with OpenType (including TrueType) fonts with the help of the `freetype2` library. Moreover, two patches can be used together.

Unfortunately, `ovp2ovf` 2.0 included in recent \TeX distributions (e.g., `te \TeX 2.x`) does not seem

⁷ Available from the FreeType project, <http://www.freetype.org>.

⁸ Available from the Korean \TeX Users group, <http://ftp.ktug.or.kr/pub/ktug/freetype/contrib/ttf2pk-1.5-20020430.patch>.

⁹ Available as http://chem.skku.ac.kr/~wkpark/project/ktug/ttf2pk-freetype2_20030314.tgz.

to work correctly, so the previous version 1.x must be used.

9 Asian Font Packs and DVIPDF x

A solution avoiding the problems mentioned above is to use the CJK fonts included in the Asian font packs of Adobe (Acrobat) Reader as non-embedded fonts when making PDF output.

It is well known that Adobe Reader can display and print several common fonts even if they are not embedded in the document. These are fourteen base Latin fonts, such as Times, Helvetica, and Courier — and several CJK fonts, if Asian font packs¹⁰ are installed. These packs have been available free of charge since the era of Adobe Acrobat Reader 4. Four are available: Chinese Simplified, Chinese Traditional, Japanese, and Korean. Moreover, Adobe Reader 6 downloads the appropriate font packs on demand when a document containing non-embedded CJK characters is opened. Note that these fonts are licensed solely for use with Adobe Readers.

Professional CJK typesetting requires at least two font families: serif and sans serif. As of Adobe Acrobat Reader 4, Asian font packs, except for Chinese Simplified, included both families, but newer packs include only a serif family. However, newer versions of Adobe Reader can automatically substitute a missing CJK font by another CJK font installed in the operating system, so displaying both families is possible on most platforms.

If the CJK fonts included in Asian font packs are to be used, there is no need to embed the fonts when making PDF output. The PDF file should contain the font names and code points only. Some ‘generic’ font names are given in Table 1, which can be handled by Acrobat Reader 4 and later. However, these names depend on the PDF viewers.¹¹ Note that the names are not necessarily true font names. For example, `Ryumin-Light` and `GothicBBB-Medium` are the names of commercial (rather expensive) Japanese fonts. They are installed in every genuine (expensive) Japanese PostScript printer. PDF readers and PostScript-compatible low-cost printers accept these names but use compatible typefaces instead.

While \TeX generates DVI output only, `pdf \TeX` generates both DVI and PDF output. But Omega and `p \TeX` do not have counterparts generating PDF

¹⁰ Asian font packs for Adobe Acrobat Reader 5.x and Adobe Reader 6.0, Windows and Unix versions, can be downloaded from <http://www.adobe.com/products/acrobat/acrrasianfontpack.html>. For Mac OS, an optional component is provided at the time of download.

¹¹ For example, these names are hard coded in the executable file of Adobe (Acrobat) Reader, and each version has different names.

Table 1: Generic CJK font names

	Serif	Sans Serif
Chinese Simplified	STSong-Light	STHeiti-Regular
Chinese Traditional	MSung-Light	MHei-Medium
Japanese	Ryumin-Light	GothicBBB-Medium
Korean	HYSMyeongJo-Medium	HYGoThic-Medium

output yet. One solution is DVIPDFMx [2], an extension of `dvipdfm`,¹² developed by Shunsaku Hirata and one of the authors, Jin-Hwan Cho.

10 Conclusion

We have shown how Omega, with CJK-ΩTP, can be used for the production of quality PDF documents using the CJK languages.

CJK-ΩTP, as it stands, is poorly tested and documented. Especially needed are examples of Chinese typesetting, in which the present authors are

¹² The utility `dvipdfm` is a DVI to PDF translator developed by Mark A. Wicks. The latest version, 0.13.2c, was released in 2001. Available from <http://gaspra.kettering.edu/dvipdfm/>.

barely qualified. In due course, we hope to upload CJK-ΩTP to CTAN.

References

- [1] ASCII Corporation. ASCII Nihongo T_EX (Publishing T_EX). <http://www.ascii.co.jp/pb/ptex/>.
- [2] Jin-Hwan Cho and Shunsaku Hirata. The DVIPDFMx Project. <http://project.ktug.or.kr/dvipdfmx/>.
- [3] John Plaice and Yannis Haralambous. The Omega Typesetting and Document Processing System. <http://omega.enstb.org>.

MIBIB \TeX : Beyond \LaTeX

Jean-Michel Hufflen

LIFC (FRE CNRS 2661)

University of Franche-Comté

16, route de Gray

25030 Besançon Cedex

France

hufflen@lifc.univ-fcomte.fr

<http://lifc.univ-fcomte.fr/~hufflen>

Abstract

This article sums up our experience with MIBIB \TeX , our multilingual implementation of BIB \TeX , and points out some possible improvements for better cooperation between \LaTeX and MIBIB \TeX . Also, MIBIB \TeX may be used to generate bibliographies written according to other formalisms, especially formalisms related to XML, and we give some ways to ease that.

KEYWORDS: Bibliographies, multilingual features, BIB \TeX , MIBIB \TeX , `bst`, `nbst`, XML, XSLT, XSL-FO, DocBook.

1 Introduction

MIBIB \TeX (for ‘MultiLingual BIB \TeX ’) is a reimplementaion of BIB \TeX [21], the bibliography processor associated with \LaTeX [19]. The project began in October 2000, and has resulted in two experimental versions [9, 11] and the present version (1.3), that will be available publicly by the time this article appears. As we explained in [15], a prototype using the Scheme programming language is working whilst we are developing a more robust program written in C. The prototype has allowed us to get some experience with real-sized bibliographies: this is the purpose of the first part of this article, after a short review of the *modus operandi* of MIBIB \TeX .

MIBIB \TeX ’s present version no longer uses the `bst` language of BIB \TeX for bibliography styles [20]. Such `.bst` files were used in MIBIB \TeX ’s first version, but since this old-fashioned language, based on simple stack manipulations, is not modular, we quickly realised that this choice would have led us to styles that were too complicated [12]. Thus, Version 1.3 uses the `nbst` (for ‘New Bibliography STyles’) language, described in [13] and similar to XSLT,¹ the language of transformations designed for XML texts [32]. More precisely, MIBIB \TeX 1.3 uses XML² as a central formalism in the sense that parsing files containing *bibliographical entries* (`.bib` files) results

in a DOM³ tree. Bibliography styles written using `nbst` are XML texts, too.

Of course, `nbst` can be used to generate bibliographies for documents other than those processed with \LaTeX .⁴ In particular, `nbst` eases the generation of bibliographies for documents written using XML-like syntax. Nevertheless, dealing with `.bib` files raises some problems: we go into them thoroughly in Section 4.

Reading this article requires only a basic knowledge of \LaTeX , BIB \TeX and XML. Some examples given in the next section will use the commands provided by the multilingual `babel` package of \LaTeX 2 ϵ [2]. Other examples given in Section 4 will use the Scheme programming language, but if need be, referring to an introductory book such as [28] is sufficient to understand them.

2 Architecture of MIBIB \TeX

2.1 How MIBIB \TeX Works

As a simple example of using MIBIB \TeX with \LaTeX , let us consider the `silke1988` bibliographical entry given in Figure 1. As we explain in [15], the sequence ‘[...] ! *idf*’ is one of the multilingual features

¹ EXtensible Stylesheet Language Transformations.

² EXtensible Markup Language. A good introduction to this formalism issued by the W3C (World Wide Web Consortium) is [24].

³ Document Object Model. This is a W3C recommendation for a standard tree-based programming approach [24, p. 306–308], very often used to implement XML trees.

⁴ This is also the case with the `bst` language of BIB \TeX , but in practice, it seems that this feature has not been used, except for documents written in SCRIBE [25], a predecessor of \LaTeX .

```
@BOOK{silke1988,
  AUTHOR = {James~R. Silke},
  TITLE = {Prisoner of the Horned Helmet},
  PUBLISHER = {Grafton Books},
  YEAR = 1988,
  NUMBER = 1,
  SERIES = {Frank Frazetta's Death Dealer},
  NOTE = {[Pas de traduction fran\c{c}aise
          connue] ! french
          [Keine deutsche Übersetzung
          ! german]},
  LANGUAGE = english}
```

Figure 1: Example of a bibliographical entry in MIBIB \TeX .

provided by MIBIB \TeX , defining a string to be included when the language of a corresponding reference, appearing within a bibliography, is *idf*. So if this entry is cited throughout a document written in French and the ‘References’ section is also written in French, it will appear as:

- [1] James R. SILKE: *Prisoner of the Horned Helmet*. N° 1 in *Frank Frazetta's Death Dealer*. Grafton Books, 1988. Pas de traduction française connue.

Here and in the bibliography of this article, we use a ‘plain’ style, that is, references are labelled with numbers. More precisely, the source processed by \LaTeX , included into the .bbl file generated by MIBIB \TeX , is:

```
\begin{thebibliography}{...}
...
\bibitem{silke1988}
\begin{otherlanguage*}{english}
James~R. \textsc{Silke}: \emph{Prisoner of the
Horned Helmet}.
\foreignlanguage{french}{\bblno~1 \bblbf}
\emph{Frank Frazetta's Death Dealer}. Grafton
Books, 1988. \foreignlanguage{french}{Pas de
traduction fran\c{c}aise connue}.
\end{otherlanguage*}
...
\end{thebibliography}
```

Let us examine this source text. We can notice the use of additional \LaTeX commands to put some keywords (‘ \bblin ’ for ‘*in*’, ‘ \bblno ’ for ‘N°’, that is, ‘number’ in French). In [14], we explain how to put them into action within \LaTeX and how MIBIB \TeX uses them. This source also shows how English words, originating from an entry in English (see the value of the LANGUAGE field in Figure 1), are processed. If the document uses the `babel` package, and if the french option of this package is selected, we use the `\foreignlanguage` command of this pack-

```
<book id="silke1988" language="english">
  <author>
    <name>
      <personname>
        <first>James&#160;R.</first>
        <last>Silke</last>
      </personname>
    </name>
  </author>
  <title>Prisoner of the Horned Helmet</title>
  <publisher>Grafton Books</publisher>
  <year>1988</year>
  <number>1</number>
  <series>
    Frank Frazetta's Death Dealer
  </series>
  <note>
    <group language="french">
      Pas de traduction française connue
    </group>
    <group language="german">
      Keine deutsche Übersetzung
    </group>
  </note>
</book>
```

Figure 2: The XML tree corresponding to the entry of Figure 1.

age [2], as shown above. Users do not have to select its `english` option; if it is not active, the source text generated by MIBIB \TeX looks like:

```
\bibitem{silke1988}James~R. \textsc{Silke}:
\emph{Prisoner of the Horned Helmet}. \bblno~1
\bblbf\ \emph{Frank Frazetta's Death Dealer}.
Grafton Books, 1988. Pas de traduction
fran\c{c}aise connue.
```

but the English words belonging to this reference will be taken as French by \LaTeX and thus may be processed or hyphenated incorrectly.

2.2 The Modules of MIBIB \TeX

As mentioned in the introduction, parsing a .bib file results in a DOM tree. In fact, .bib files are processed as if they were XML trees, but without whitespace nodes.⁵ Following this approach, the entry `silke1988` given in Figure 1 is viewed as the tree of Figure 2, except that the whitespace nodes that an XML parser would produce are excluded.

⁵ These are text nodes whose contents are only whitespace characters, originating from what has been typed between two tags [27, p. 25–26]. For example, if the XML text of Figure 2 is parsed, there is a whitespace node, containing a newline and four space characters between the opening tags `<author>` and `<name>`. XML parsers are bound by the ‘all text counts’ constraint included in the XML specification [33, § 2.10], and cannot ignore such whitespace characters.

We can see that some L_{ATEX} commands and special characters are converted according to the conventions of XML.

- The commands used for accents and special letters are replaced by the letter itself. This poses no problem since DOM trees are encoded in Unicode [29]. As an example, the ‘\c{c}’ sequence in the value of the NOTE field in Figure 1 is replaced by ‘ç’ in Figure 2. (By the way, let us remark that MIBIB_{TEX} can handle the 8-bit latin1 encoding:⁶ notice the ‘Û’ character inside this value.)
- Likewise, the commands:
 - ‘\ ’ for a simple space character,
 - ‘\ ’ for an end-of-line character,

and the sequences of characters:

- ‘~’, for an unbreakable space character,
- ‘--’, and ‘---’ for dash characters,

are replaced by the corresponding Unicode values for these characters:⁷

```
&#32; &#10; &#160; &#8211; &#8212;
```

An example is given by the value of the AUTHOR field, see Figures 1 & 2.

- Some characters escaped in L_{ATEX} (for example, ‘\$’, ‘%’, ‘&’) lose the escape character:

```
\% ==> %
```

The escape is restored if MIBIB_{TEX} generates a .bbl file to be processed by L_{ATEX}. Other characters are replaced by a reference to a character entity:⁸

```
\& ==> &amp; < ==> &lt; > ==> &gt;
```

- Balanced delimiters for quotations (“ ” and “ ’ ” or “ ‘ ” and “ ’ ’) are replaced by an **emph** element:⁹

```
‘Tooth and Claw’ ==>
<emph emf=’no’ quotedf=’yes’>
  Tooth and Claw
</emph>
```

If ‘ ’ or “ ” characters are unbalanced, they are replaced by references to character entities used in XML documents:

```
’ ==> &apos; " ==> &quot;
```

⁶ See [7, Table C.4] for more details.

⁷ That was not the case in earlier versions; for instance, [12, Figure 3] improperly includes a tilde character in a text node. This bug was fixed at the end of 2003.

⁸ See [24, p. 48] for more details.

⁹ ‘**emph**’ is of course for ‘emphasise’: all the attributes (for example, ‘**quotedf**’ for ‘quoted-flag’, used for specifying a quotation) default to **no**, except **emf**, which defaults to **yes**. The complete specification is given in [10].

Such an XML tree, resulting from our parser, may be validated using a DTD;¹⁰ more precisely, by a revised version of the DTD sketched in [10].

Some examples of using **nbst** for bibliography styles are given in [12, 13, 14]. We give another example in Figure 3. We can see that this language is close to XSLT and it uses *path expressions* as in the XPath language [31]. Also, the example shows how multilingual features (for example, the sequence ‘[...] ! ...’) are processed: we use some external functions in order to determine which L_{ATEX} command can be used to switch to another language. These external functions are written using the language of MIBIB_{TEX}’s implementation: Scheme for the prototype, C for the final program.

3 MIBIB_{TEX} with L_{ATEX}

When BIB_{TEX} generates a .bbl file, it does not use the source file processed by L_{ATEX}, but only the auxiliary (.aux) file, in which the definition of all the labels provided by the commands `\label` and `\bibitem` is stored. This file also contains the name of the bibliography style to be used and the paths of bibliography data bases to be searched, so BIB_{TEX} need not look at any other file.

This is not true for MIBIB_{TEX}. It still uses the .aux file as far as possible, but it also has to determine which multilingual packages are used: first of all **babel**, but also some packages devoted to particular languages: **french** [6], **german** [23], ... So we have to do a *partial parsing* of the .tex file for that. For better co-operation between L_{ATEX} and MIBIB_{TEX}, this could be improved, in that information about multilingual packages used, and languages available, could be put in the .aux file. In fact, the external functions of our new bibliography styles are only used to manage information extracted from a .tex file. Expressing such operations using **nbst** would be tedious.

Another improvement regarding the natural languages known by L_{ATEX} would be a connection between:

- the language codes used in XML, specified by means of a two-letter language abbreviation, optionally followed by a two-letter country code [1] (for example, ‘**de**’ for ‘*deutsch*’ (‘German’), ‘**en-UK**’, ‘**en-US**, etc.’); and
- the *resources* usable to write texts in these languages.

For example, a default framework could be the use of the **babel** package, and ‘**de**’ would get access to

¹⁰ **Document Type Definition**. A DTD defines a document markup model [24, Ch.5].

```

<nbst:template match="group">
  <nbst:if test="@language=$language">
    <!-- The $language variable is set to the current language. -->
    <nbst:value-of select="call(language_open_change,@language)"/>
    <!-- If the babel package is used and a known option has been selected, this external function
         writes the \foreignlanguage command...
    -->
    <nbst:apply-templates use-language="@language"/>
    <nbst:value-of select="call(language_close_change,@language)"/>
    <!-- ... and this external function puts a closing brace. -->
  </nbst:if>
</nbst:template>

```

Figure 3: Example of calling an external function.

the german option of this package, although it could be redefined to use the *ad hoc* package name `german`. In the future, such a framework would allow us to homogenise all the notations for natural languages to those of XML. In addition, let us notice that ConTeXt¹¹ [8], already uses these two-letter codes in its `\selectlanguage` command.

And last but not least, auxiliary files should include information about the encoding used in the source text. As can be seen in the examples of Section 2.1, accented letters are replaced by the commands used to produce them in L^AT_EX, even though L^AT_EX can of course handle 8-bit encodings (provided that the `inputenc` package is loaded with the right option). This is to avoid encoding problems. In addition, such information would ease the processing of languages written using non-Latin alphabets.

4 Towards the XML World

Since a `.bib` file can be processed as an XML tree by a bibliography style written in `nbst`, MIBIB_TE_X opens a window on XML's world. A converter from `.bib` files to a file written using HTML,¹² the language of Web pages, becomes easy to write. So does a tool to write a bibliography as an XSL-FO¹³ document [34]. More precisely, we give in Figure 4 an example of using the root element of `nbst`. Possible values for the `method` of the `nbst:output` element are:

```

<nbst:bst version="1.3" id="plain" xmlns:nbst=
"http://lifc.univ-fcomte.fr/~hufflen/mlbibtex"
>
  <nbst:output method="LaTeX"/>
  ...
</nbst:bst>

```

Figure 4: Root element for a bibliography style written using `nbst`.

LaTeX xml html text

Nevertheless, this approach has an important limitation in practice. Since BIB_TE_X has traditionally been used to generate files suitable for L^AT_EX, users often put L^AT_EX commands inside values of BIB_TE_X fields.¹⁴ For example:

```
ORGANIZATION = {\textsc{tug}}
```

In such a case, we would have to write a mini-L^AT_EX program (or perhaps a new output mode for L^AT_EX) that would transform such a value into a string suitable for an XML parser.

The problem is more complicated when commands are defined by end-users. For instance:

```
ORGANIZATION = {\logo{tug}}
```

works with BIB_TE_X—or MIBIB_TE_X when we use it for generating L^AT_EX output—even though `\logo` has an arbitrary definition; for example,

```
\newcommand{\logo}[1]{\textsc{#1}}
```

according to L^AT_EX's conventions, or:

```
\def\logo#1{\textsc{#1}}
```

if a style close to plain T_EX is used. Likewise, such commands can be known when an output file from MIBIB_TE_X is processed by ConTeXt.

¹¹ T_EX, defined by Donald E. Knuth [18], provides a general framework to format texts. To be fit for use, the definitions of this framework need to be organised in a *format*. Two such formats are plain T_EX and L^AT_EX, and another is ConTeXt, created by Hans Hagen.

¹² HyperText Markup Language.

¹³ EXtensible Stylesheet Language—Formatting Objects: this language aims to describe high-quality print outputs. Such documents can be processed by the shell command `xmltex` (resp. the shell command `pdfxmltex`) from PassiveT_EX [22, p. 180] to get `.dvi` files (resp. `.pdf` files).

¹⁴ The author personally confesses to using many `\foreignlanguage` commands within the values of BIB_TE_X fields, before deciding to develop MIBIB_TE_X.

```

<bibliography>
  <title>References</title>
  <biblioentry>
    <abbrev>silke1989</abbrev>
    <authorgroup>
      <author>
        <firstname>James R.</firstname>
        <surname>Silke</surname>
      </author>
    </authorgroup>
    <copyright><year>1989</year></copyright>
    <isbn>0-586-07018-4</isbn>
    <publisher>
      <publishername>
        Grafton Books
      </publishername>
    </publisher>
    <title>Lords of Destruction</title>
    <seriesinfo>
      <title>
        <othercredit>
          <firstname>Frank</firstname>
          <surname>Frazetta</surname>
        </othercredit>'s Death Dealer
      </title>
      <volumenum>2</volumenum>
    </seriesinfo>
  </biblioentry>
</bibliography>

```

Figure 5: The bibliographical reference from Figure 1 expressed in DocBook. Note the *ad hoc* tag `<othercredit>`.

Moreover, let us consider the bibliographical reference given in Figure 5, according to the conventions of DocBook, a system for writing structured documents [36] (we use the conventions of the XML version of DocBook, described in [26]). We can see that some information is more precise than that provided in Figure 1. But there are still complexities: the person name given in the value of the `SERIES` field is surrounded by an *ad hoc* element in the DocBook version.

If we want to take advantage of the expressive power of DocBook, we can:

- directly process an XML file for bibliographical entries. In this case, our DTD should be extended; that is possible, but we still need a solution to process the huge number of existing `.bib` files;
- introduce some new syntax inside `.bib` files, that might be complicated and thus perhaps unused in practice,

- introduce new L^AT_EX commands, to process like the `\logo` example mentioned above.

We have experimentally gone quite far in the third direction, which also allows to us to deal with the L^AT_EX commands already in `.bib` files. In Figure 6, we give some examples of such processing, as implemented in the prototype.¹⁵

As can be seen, we have defined a new function in Scheme, named `define-pattern`, with two formal arguments. The first is a string viewed as a *pattern*, following the conventions of T_EX for defining commands, that is, the arguments of a command are denoted by ‘#1’, ‘#2’, ... (cf. [18, Ch. 20]). The second argument may also be a string, in which case it specifies a replacement. The arguments of the corresponding command are processed recursively. In case of conflict among patterns, the longest is chosen. So, the pattern `"\\logo{#1}"`¹⁶ takes precedence over the pattern `"{#1}"`.

If the second argument of the `define-pattern` function is not a string, it must be a zero-argument function that results in a string. In this case, all the operations must be specified explicitly, using the following functions we wrote:

`pattern-matches?` returns a *true* value if its first argument matches the second, a *false* value otherwise;

`pattern-process` recursively processes its only argument, after replacing sub-patterns by corresponding values;¹⁷

`pattern-replace` replaces the sub-patterns of its argument by corresponding values; these value are not processed, just replaced *verbatim*.

Whether given directly as the second argument to `define-pattern` or resulting from applying a zero-argument function, the string must be well-formed w.r.t. XML’s conventions, that is, tags must be balanced, attributes must be well-formed, etc. In other words, such a string must be acceptable to an XML parser: in our case, the parser is SSAX¹⁸ [17].

The examples given in Figure 6 allow us to see that we can deal with simple commands, like:

```
\logo{...} ==> <emph ...>...</emph>
```

¹⁵ This feature has not yet been implemented in the final version.

¹⁶ Let us recall that in Scheme, the backslash character (`\`) is used to escape special characters in string constants. To include it within a string, it must itself be escaped.

¹⁷ In fact, using a string *s* as a second argument of `define-pattern` yields the evaluation of the expression `(lambda () (pattern-process s))`.

¹⁸ Scheme implementation of SAX. ‘SAX’ is for ‘Simple API (Application Programming Interface) for XML’: this name denotes a kind of parser, see [24, p. 290–292].

```

(define-pattern "{#1}"
  ;; The asitis element is used for words that should never be uncapitalised, that is, proper names. In BIBTEX,
  ;; we specify such behaviour by surrounding words by additional braces.
  "<asitis>#1</asitis>")

(define-pattern "\\logo{#1}" "<emph emf='no' scf='yes'>#1</emph>")

(define-pattern "\\foreignlanguage{#1}{#2}"
  "<foreigngroup language='#1'>#2</foreigngroup>")

(define-pattern "\\iflanguage{#1}{#2}{#3}"
  (lambda () ; Zero-argument function.
    (string-append ; Concatenation of strings.
      "<nonemptyinformation>"
      "<group language='" (pattern-replace "#1") "'>" (pattern-process "#2")
      "</group>"
      (let loop ((pattern (pattern-replace "#3")))
        ;; This form—named let (cf. [28, Exercise 14.8])—defines an internal function loop and
        ;; launches its first call:
        (if (pattern-matches? "\\iflanguage{#4}{#5}{#6}" pattern)
            (string-append "<group language='" (pattern-replace "#4") "'>"
                           (pattern-process "#5")
                           "</group>"
                           ;; The internal function loop is recursively called with a new value:
                           (loop (pattern-replace "#6")))
            (string-append "<group>" (pattern-process pattern) "</group>")))
      "</nonemptyinformation>")))

```

Figure 6: Patterns for some L^AT_EX commands in Scheme.

as well as more complicated cases, like a cascade of `\iflanguage` commands [2]:

```

\iflanguage{...}{...}{%
  \iflanguage{...}{...}{ ... }}

```

which becomes:

```

<nonemptyinformation>
  <group language='...'>...</group>
  <group language='...'>...</group>
  ...
</nonemptyinformation>

```

The `nonemptyinformation` element is used for information that must be output, possibly in a default language if no translation into the current language is available.

What we do by means of our `define-pattern` function is like the additional procedures in Perl¹⁹ that the converter `LaTeX2HTML` [4] can use to translate additional commands.

5 Conclusion

Managing several formalisms can be tedious. This fact was one of main elements in XML's design: giving a central formalism, able to be used for representing trees, and allowing many tools using different formalisms to communicate.

¹⁹ Practical Extraction and Report Language.

BIB_TE_X deals with three formalisms: `.aux` files, `.bib` files and `.bst` files. As Jonathan Fine notes in [5], the applications devoted to a particular formalism cannot be shared with other applications. MIBIB_TE_X attempts to use XML as far as possible, although there is still much to do. For example, defining a syntax for the entries for which we are looking, when using MIBIB_TE_X to generate XSL-FO or DocBook documents. (For our tests, this list of entry names is simply given on the command line).

The next step will probably be a more intensive use of XML, that is, the direct writing of bibliographical entries using XML conventions. For this, we need something more powerful than DTDs, with a richer type structure, namely, *schemas*.²⁰ In addition, we should be able to easily add new fields to bibliographical entries: the example given using DocBook shows that additional information must be able to be supplied to take advantage of the expressive power of this system. But such additions are

²⁰ Schemas have more expressive power than DTDs, because they allow users to define types precisely, which in turn makes for a better validation of an XML text. In addition, this approach is more homogeneous since schemas are XML texts, whereas DTDs are not.

There are currently four ways to specify schemas: Relax NG [3], Schematron [16], Examplotron [30], XML Schema [35]. At present, it seems to us that XML Schema is the most suitable for describing bibliographical entries.

difficult to model with DTDs.²¹ We are presently going thoroughly into replacing our DTD by a schema; when this work reaches maturity, bibliographical entries using XML syntax could be directly validated using schemas.

6 Acknowledgements

Many thanks to Karl Berry for his patience while waiting for this article. In addition he proofread a first version and gave me many constructive criticisms. Thanks also to Barbara Beeton.

References

- [1] Harald Tveit ALVSTRAND: *Request for Comments: 1766. Tags for the Identification of Languages*. UNINETT, Network Working Group. March 1995. <http://www.cis.ohio-state.edu/cgi-bin/rfc/rfc1766.html>.
- [2] Joannes BRAAMS: *Babel, a Multilingual Package for Use with \LaTeX 's Standard Document Classes*. Version 3.7. May 2002. CTAN:macros/latex/required/babel/babel.dvi.
- [3] James CLARK *et al.*: *Relax NG*. <http://www.oasis-open.org/committees/relax-ng/>. 2002.
- [4] Nicos DRAKOS: *The \LaTeX 2HTML Translator*. March 1999. Computer Based Learning Unit, University of Leeds.
- [5] Jonathan FINE: “ \TeX as a Callable Function”. In: *Euro \TeX 2002*, (pp. 26–30). Bachotek, Poland. April 2002.
- [6] Bernard GAULLE : *Notice d'utilisation du style french multilingue pour \LaTeX* . Version pro V5.01. Janvier 2001. CTAN:loria/language/french/pro/french/ALIRE.pdf.
- [7] Michel GOOSSENS, Sebastian RAHTZ and Frank MITTELBACH: *The \LaTeX Graphics Companion. Illustrating Documents with \TeX and PostScript*. Addison-Wesley Publishing Company, Reading, Massachusetts. March 1997.
- [8] Hans HAGEN: *Con \TeX t, the Manual*. November 2001. <http://www.pragma-ade.com>.
- [9] Jean-Michel HUFFLEN: “MIBIB \TeX : A New Implementation of BIB \TeX ”. In: *Euro \TeX 2001*, (pp. 74–94). Kerkrade, The Netherlands. September 2001.
- [10] Jean-Michel HUFFLEN: “Multilingual Features for Bibliography Programs: From XML to MIBIB \TeX ”. In: *Euro \TeX 2002*, (pp. 46–59). Bachotek, Poland. April 2002.
- [11] Jean-Michel HUFFLEN: *Towards MIBIB \TeX 's Versions 1.2 & 1.3*. Ma \TeX Conference. Budapest, Hungary. November 2002.
- [12] Jean-Michel HUFFLEN: “European Bibliography Styles and MIBIB \TeX ”. Euro \TeX 2003, Brest, France. June 2003. (To appear in *TUGboat*.)
- [13] Jean-Michel HUFFLEN: *MIBIB \TeX 's Version 1.3*. TUG 2003, Outrigger Waikoloa Beach Resort, Hawaii. July 2003.
- [14] Jean-Michel HUFFLEN: “Making MIBIB \TeX Fit for a Particular Language. Example of the Polish Language”. *Biuletyn GUST*. Forthcoming. Presented at the Bacho \TeX 2003 conference. 2004.
- [15] Jean-Michel HUFFLEN: “A Tour around MIBIB \TeX and Its Implementation(s)”. *Biuletyn GUST*, Vol. 20, pp. 21–28. In *Proc. Bacho \TeX Conference*. April 2004.
- [16] ISO/IEC 19757: *The Schematron. An XML Structure Validation Language Using Patterns in Trees*. <http://www.ascc.net/xml/resource/schematron/schematron.html>. June 2003.
- [17] Oleg KISELYOV: “A Better XML Parser through Functional Programming”. In: *4th International Symposium on Practical Aspects of Declarative Languages*, Vol. 2257 of *Lecture Notes in Computer Science*. Springer-Verlag. 2002.
- [18] Donald Ervin KNUTH: *Computers & Typesetting. Vol. A: The \TeX book*. Addison-Wesley Publishing Company, Reading, Massachusetts. 1984.
- [19] Leslie LAMPORT: *\LaTeX . A Document Preparation System. User's Guide and Reference Manual*. Addison-Wesley Publishing Company, Reading, Massachusetts. 1994.
- [20] Oren PATASHNIK: *Designing BIB \TeX styles*. February 1988. Part of the BIB \TeX distribution.
- [21] Oren PATASHNIK: *BIB \TeX ing*. February 1988. Part of the BIB \TeX distribution.
- [22] Dave PAWSON: *XSL-FO*. O'Reilly & Associates, Inc. August 2002.
- [23] Bernd RAICHLE: *Die Makropakete „german“ und „ngerman“ für \LaTeX 2 ϵ , \LaTeX 2.09, Plain- \TeX and andere darauf Basierende Formate*. Version 2.5. Juli 1998. Im Software \LaTeX .
- [24] Erik T. RAY: *Learning XML*. O'Reilly & Associates, Inc. January 2001.

²¹ Whereas that is easy with ‘old’ BIB \TeX , provided that you use a bibliography style able to deal with additional fields.

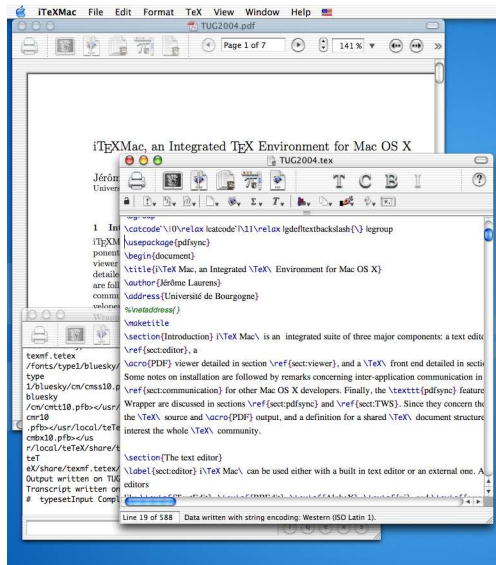
- [25] Brian Keith REID: *SCRIBE Document Production System User Manual*. Technical Report, Unilogic, Ltd. 1984.
- [26] Thomas SCHRAITL: *DocBook-XML—Medien-neutrales und plattformunabhängiges Publizieren*. SuSE Press. 2004.
- [27] John E. SIMPSON: *XPath and XPointer*. O'Reilly & Associates, Inc. August 2002.
- [28] George SPRINGER and Daniel P. FRIEDMAN: *Scheme and the Art of Programming*. The MIT Press, McGraw-Hill Book Company. 1989.
- [29] *The Unicode Standard Version 3.0*. Addison-Wesley. February 2000.
- [30] Eric VAN DER VLIST: *Examplotron*. <http://examplotron.org>. February 2003.
- [31] W3C: *XML Path Language (XPath). Version 1.0*. W3C Recommendation. Edited by James Clark and Steve DeRose. November 1999. <http://www.w3.org/TR/1999/REC-xpath-19991116>.
- [32] W3C: *XSL Transformations (XSLT). Version 1.0*. W3C Recommendation. Written by Sharon Adler, Anders Berglund, Jeff Caruso, Stephen Deach, Tony Graham, Paul Grosso, Eduardo Gutentag, Alex Milowski, Scott Parnell, Jeremy Richman and Steve Zilles. November 1999. <http://www.w3.org/TR/1999/REC-xslt-19991116>.
- [33] W3C: *Extensible Markup Language (XML) 1.0 (Second Edition)*. W3C Recommendation. Edited by Tim Bray, Jean Paoli, C. M. Sperberg-McQueen and Eve Maler. October 2000. <http://www.w3.org/TR/2000/REC-xml-20001006>.
- [34] W3C: *Extensible Stylesheet Language (XSL). Version 1.0*. W3C Recommendation. Edited by James Clark. October 2001. <http://www.w3.org/TR/2001/REC-xsl-20011015/>.
- [35] W3C: *XML Schema*. November 2003. <http://www.w3.org/XML/Schema>.
- [36] Norman WALSH and Leonard MUELLNER: *DocBook: The Definitive Guide*. O'Reilly & Associates, Inc. October 1999.

iTeXMac: An Integrated TeX Environment for Mac OS X

Jérôme Laurens

Université de Bourgogne

jerome.laurens@u-bourgogne.fr



1 Introduction

iTeXMac is an integrated suite of three major components: a text editor detailed in section 2, a PDF viewer detailed in section 3, and a TeX front end detailed in section 4. Some notes on installation are followed by remarks concerning inter-application communication in section 6 for other Mac OS X developers. Finally, the `pdfsync` feature and the TeX Wrapper are discussed in sections 7 and 8. Since they concern the synchronization between the TeX source and PDF output, and a definition for a shared TeX document structure, both will certainly interest the whole TeX community.

2 The Text Editor

iTeXMac can be used either with a built-in text editor or an external one. All standard text editors like TextEdit, BBEdit, AlphaX, vi, and emacs are supported and configuring iTeXMac for other editors is very easy, even when coming from the X11 world.

The built-in text editor comes with flavours similar to emacs and AlphaX modes. It relies on a plug-in architecture that allows very different kinds of user interfaces according to the type of the file being edited. Whereas AlphaX uses Tcl and emacs uses Lisp, iTeXMac utilizes the benefits of Objective-C

bundles, giving plug-ins great potential power together with the application.

Among the standard features shared by advanced text editors (like key binding management, advanced regular expressions, command completion), an interesting feature of iTeXMac's text editor is the syntax parsing policy. The syntax highlighting deeply depends on the kind of text edited, whether it is Plain, L^AT_EX or METAPOST (support for HTML is planned). The text properties used for highlighting include not only the color of the text, but also the font, the background color and some formatting properties.

Moreover, the command shortcuts that refer to mathematical and text symbols are replaced by the glyph they represent, thus replacing `\alpha` with α and so on. Conversely the built-in editor can show a character palette with 42 menus gathering text and mathematical symbols, as they would appear in the output. The editor thus serves as a graphical front-end to the standard L^AT_EX packages, `amsfonts.sty`, `amssymb.sty`, `mathbb.sty`, `mathrsfs.sty`, `marvosym.sty` and `wasysym.sty`, which makes thousands of symbols available with just one click. The result is a text editor that contains much more WYSIWYG than others, with no source file format requirement.

There is also advanced management of string encoding, and iTeXMac supports more than 80 of them with an efficient user interface. The text files are scanned for hints about the text encoding:

L ^A T _E X	<code>\usepackage[<i>encoding</i>]{inputenc}</code>
ConT _E Xt	<code>\enableregime[<i>encoding</i>]</code>
emacs	<code>%-*-coding:character encoding;-*-</code>
	<code>%!iTeXMac(charset):</code> <i>character encoding</i>
	Mac OS X hidden internals

But this is not user friendly practice and will be enhanced by the forthcoming discussion of TeX wrappers in section 8.

Spell checking for TeX input is available with Rick Zaccane's L^AT_EX aware Excalibur¹ and Anton

¹ <http://www.eg.bucknell.edu/~excalibr/>

Leuski's \TeX aware `cocoAspell`², a port of the Free and Open Source spell checker `aspell`. The latter also knows about HTML code and is integrated to Mac OS X allowing `i \TeX Mac` to check spelling as you type, with misspelled words being underlined in red. While this service is provided to all applications for free, `i \TeX Mac` is the only one that truly enables the \TeX support by managing the language and the list of known words on a file by file basis using \TeX Wrappers.

3 The PDF Viewer

`i \TeX Mac` can be used either with a built-in PDF viewer or an external one. The built-in viewer lacks many advanced features of Acrobat or Preview, but it updates the display automatically when a PDF file has been changed externally. Moreover, it allows you to make selections and export them to other applications, like Word, TextEdit or Keynote for example. Finally, it has support for the useful PDF synchronization discussed below and is very well integrated in the suite.

`i \TeX Mac` can open PS, EPS and DVI files with a double click, by first converting them to PDF. It thus plays the role of a POSTSCRIPT or a DVI viewer. This feature is now partially obsolete since Mac OS X version 10.3 provides its own PS to PDF translator used by the Preview application shipped with the system.

4 The \TeX Front End

This component of the software serves two different purposes. On one hand it is a bridge between the user and the utilities of a standard \TeX distribution: a graphical user interface for the commands `tex`, `latex`, `pdftex`, and so on. On the other hand, it has to properly manage the different kinds of documents one wants to typeset.

Actually, the `i \TeX Mac` interface with its underlying \TeX distribution is fairly simple. Five basic actions are connected to menu items, toolbar buttons or command shortcuts to

- typeset (e.g., running `latex` once, or twice in advanced mode)
- make the bibliography (e.g. running `bibtex`)
- make the index (e.g., running `makeindex`)
- render graphics (e.g., running `dvipdf`)

All these actions are connected to shell scripts stored on a per document basis. If necessary, the user can customize them or even change the whole process by inserting an in-line instruction at the very beginning

² <http://cocoAspell.1euski.net/>

of a source file. For example, the following directive, if present, will run `pdflatex` in escape mode.

```
%!i $\TeX$ Mac(typeset): pdflatex
--shell-escape $i $\TeX$ Input
```

The `makeindex` and `bibtex` command options can be set from panels, and other commands are supported. Moreover, the various log files are parsed, warnings and errors are highlighted with different colors and HTML links point to lines where an error occurred. Some navigation facilities from log file to output are also provided, a string like `[a number...]`, pointing to the output page.

As for documents, `i \TeX Mac` manages a list of default settings that fit a wide range of situations, including for example

- \LaTeX documents with DVI, PS or PDF engines
- books
- METAPOST documents
- Con \TeX t documents
- HTML documents
- B. Gaulle's French Pro documents.

Users can extend this list with a built-in editor, adding support for Music \TeX , maybe `gcc`, and so on.

5 Installing \TeX and `i \TeX Mac`

The starting point for a detailed documentation is the MacOS X \TeX / \LaTeX Web Site³ where one will find an overview of the \TeX related tools available on Mac OS X. As a graphical front-end, `i \TeX Mac` needs a \TeX distribution to be fully functional. Gerben Wierda maintains on his site⁴ the \TeX Live⁵ distribution and a set of useful packages. Other `te \TeX` 2.0.2 ports are available from `fink`⁶ (and through one of its graphical user interfaces, such as `finkcommander`⁷) and from Darwin Ports⁸, through a CVS interface.

The official web site of `i \TeX Mac` is hosted by the Open Source software development SourceForge website at:

<http://itexmac.sourceforge.net/>

One can find in the download section the disk images for the following products:

- `i \TeX Mac`, both stable and developer release

³ <http://www.esm.psu.edu/mac-tex/>

⁴ <http://www.rna.nl/tex.html>

⁵ <http://www.tug.org/texlive/>

⁶ <http://fink.sourceforge.net/pdb/package.php/tetex>

`tetex`

⁷ <http://finkcommander.sourceforge.net/>

⁸ <http://darwinports.opendarwin.org/ports/?by=cat&substr=print>

- an external editor for iTeXMac key binding
- the Hypertext Help With L^AT_EX wrapped as a searchable Mac OS X help file.
- the TeX Catalog On line wrapped as a searchable Mac OS X help file.
- the French L^AT_EX FAQ wrapped as a searchable Mac OS X help file.

An updater allows you to check easily for new versions. To install iTeXMac, just download the latest disk image archive, double click and follow the instructions in the read-me file.

Due to its Unix core, Mac OS X is no longer focused on only one user. To support multiple users, iTeXMac configuration files can be placed in different locations to change defaults for all, or just certain users. The search path is:

- the built-in domain as shipped with the application (with default, read-only settings)
- the network domain (`/Network/Library/Application Support/iTeXMac`), where an administrator can put material to override or augment the default behaviour of all the machines on a network
- the local domain (`/Library/Application Support/iTeXMac`), where an administrator can put material to override or augment the network or default behaviour
- the user domain (`~/Library/Application Support/iTeXMac`), where the user can put material to override or augment the local or default behaviour

This is a way to apply modifications to iTeXMac as a whole.

6 Inter-application Communication

This section describes how iTeXMac communicates with other components, in the hope that this syntax will also be used by other applications when relevant, to avoid the current situation where there are as many AppleScript syntaxes as there are available applications for TeX on the Macintosh. It also shows why iTeXMac integrates so well with other editors or viewers.

6.1 Shell Commands

iTeXMac acts as a server, such that other applications can send it messages. Each time it starts, iTeXMac installs an alias to its own binary code in `~/Library/TeX/bin/iTeXMac`. With the following syntax,⁹ either from the command line, an AppleScript or shell script, one can edit a text file at the

⁹ These commands should be entered all on one line. They are broken here due to the narrow *TUGboat* columns.

location corresponding to the given line and column numbers

```
~/Library/TeX/bin/iTeXMac edit -file
  "filename"
  -line lineNumber -column colNumber
```

The following syntax is used to display a PDF file at the location corresponding to the given line column and source file name

```
~/Library/TeX/bin/iTeXMac display -file
  "filename.pdf"
  -source "sourcename.tex"
  -line lineNumber -column colNumber
```

6.2 AppleScript

The same feature is implemented using this scripting language. It would be great for the user if other TeX front ends on Mac OS X would implement the same syntax.

```
tell application "iTeXMac" to edit
  "filename.tex"
  at line lineNumber column colNumber

tell application "iTeXMac" to display
  "filename.pdf"
  at line lineNumber column colNumber
  in source "Posix source name.tex"
```

iTeXMac support for AppleScript likewise covers the compile, bibliography and index actions. They are not given here since there is no Apple Events Suite dedicated to TeX. However, configuration files and instructions are given to let third-party applications like Alpha X, BBE_{edit} or emacs control iTeXMac using those scripts.

6.3 HTML

iTeXMac implements support for a URL scheme named `file-special` for editing, updating or displaying files, for example

```
file-special://localhost/"filename.tex";
  action=edit;line=lineNumber;
  column=columnNumber

file-special://localhost/"filename.pdf";
  action=display;line=lineNumber;
  column=columnNumber;
  source="Posix source name.tex"
```

will ask iTeXMac to edit a TeX source file or display the given file (assumed to be PDF) and when synchronization information is available, scroll to the location corresponding to the given line and column in source (assumed to be TeX). This allows adding dynamic links in HTML pages, in a TeX tutorial for example.

7 The pdfsync Feature

7.1 About Synchronization

As the \TeX typesetting system heavily relies on a page description language, there is no straightforward correspondence between a part of the output and the original description code in the input. A workaround was introduced a long time ago by commercial \TeX frontends Visual \TeX ¹⁰ and \TeX tures¹¹ with a very efficient implementation. Then \LaTeX users could access the same features—though in a less-efficient implementation—through the use of `src1tx.sty`, which added source specials in the DVI files. The command line option `-src-specials` now gives this task to the \TeX typesetting engine.

When used with an external DVI viewer or an external text editor, through an X11 server or not, $i\TeX$ Mac fully supports this kind of synchronization feature.

For the PDF file format, Piero d’Ancona and the author elaborated a strategy that works rather well for Plain \TeX , \ConTeX t and \LaTeX users. While typesetting a `foo.tex` file with \LaTeX for example, the `pdfsync` package writes extra geometry information in an auxiliary file named `foo.pdfsync`, subsequently used by the front ends to link line numbers in source documents with locations in pages of output PDF documents. $i\TeX$ Mac and \TeX Shop¹² both support `pdfsync`.

The official `pdfsync` web site is:

<http://itexmac.sourceforge.net/pdfsync.html>

It was more convenient to use an auxiliary file than to embed the geometric information in the PDF output using `pdfTeX` primitives. The output file is not polluted with extraneous information and the front ends need not parse the PDF output to retrieve such metrics.

7.2 The pdfsync Mechanism

A macro is defined to put `pdfsync` anchors at specific locations (for `hbox`’s, paragraphs and maths). There are essentially three problems we must solve: the position of an object in the PDF page is not known until the whole page is composed, the objects don’t appear linearly in the output¹³ and finally, an input file can be entirely parsed long before its contents are shipped out. To solve these, at each `pdfsync` anchor the known information (line number and source file name) is immediately written to the `pdfsync` auxiliary file and the unknown infor-

mation (location and page number) will be written only at the next ship out.

7.3 The .pdfsync File Specifications

This is an ASCII text file organized into lines. There is no required end of line marker format from among the standard ones used by operating systems.

Only the two first lines described in table 1 are required, the other ones are optional. The remaining lines are described according to their starting characters, they consist of 2 interlaced streams. A synchronous one detailed in table 2 is obtained with `\immediate\writes` and concerns the input information. An asynchronous one detailed in table 3 is obtained with delayed `\writes` and concerns the output information.

The correspondence between the two kinds of information is made through a record counter, which establishes a many-to-many mapping from line numbers in \TeX sources to positions in PDF output.

7.4 Known Problems

Unfortunately, the various `pdfsync` files for Plain, \LaTeX or \ConTeX t are not completely safe. Some compatibility problems with existing macro packages may occur. Moreover, sometimes `pdfsync` actually influences the final layout; in a case like that, it should only be used in the document preparation stage.

Another mechanism widely used by \ConTeX t makes `pdfsync` sometimes inefficient, where the macro expansion only occurs long after it has been parsed, such that the `\inputlineno` is no longer relevant and the significant line number is no longer accessible. This makes a second argument for the implementation of the `pdfsync` feature at a very low level, most certainly inside the `pdfTeX` engine itself.

8 TWS: A \TeX Wrapper Structure

In general, working with \TeX seems difficult due to the numerous auxiliary files created. Moreover, sharing \TeX documents is often delicate as soon as we do not use very standard \LaTeX . The purpose of this section is to lay the foundation for the \TeX Wrapper Structure, which aims to help the user solve these problems.

First, it is very natural to gather all the files related to one \TeX document in one folder we call a \TeX Wrapper. The file extension for this directory is `texd`, in reference to the `rtf` and `rtfd` file extensions already existing on Mac OS X. The contents of a \TeX wrapper named `document.texd` is divided according to different criteria:

¹⁰ <http://www.micropress-inc.com/>

¹¹ <http://www.bluesky.com/>

¹² <http://www.uoregon.edu/~koch/texshop>

¹³ The footnotes objects provide a good example.

Line	Format	Description	Comment
1 st	<i>jobName</i>	<i>jobName</i> : case sensitive T _E X file name	In general, the extensionless name of the file as the result of an <code>\immediate\write\file{\jobname}</code>
2 nd	version <i>V</i>	<i>V</i> : a 0 based non-negative integer	The current version is 0

Table 1: pdfsync required lines format

Line	Format	Description	Comment
“b”	b <i>name</i>	<i>name</i> : T _E X file name	T _E X is about to begin parsing <i>name</i> , all subsequent line and column numbers will refer to <i>name</i> . The path is relative to the directory containing the .pdfsync file. Path separators are the Unix “/”. The file extension is not required, “tex” is the default if necessary. Case sensitive.
“e”	e		The end of the input file has been reached. Subsequent line and column numbers now refer to the calling file. Optional, but must match a corresponding “b” line.
“l”	l <i>R L</i> l <i>R L C</i>	<i>R</i> : record number, <i>L</i> : line number, <i>C</i> : optional column number.	

Table 2: pdfsync line specifications of the synchronous stream

Line	Format	Description	Comment
“s”	s <i>S</i>	<i>S</i> : physical page number	T _E X is going to ship out a new page.
“p” “p*”	p <i>R x y</i> p* <i>R x y</i>	<i>R</i> : record number, <i>x</i> : horizontal coordinate, <i>y</i> : vertical coordinate.	Both coordinates are respectively given by <code>\the\pdfastxpos</code> and <code>\the\pdfastypos</code>

Table 3: pdfsync line specification of the asynchronous stream

- required data (source, graphics, bibliography database)
- helpful data and hints (tex, bibtex, makeindex options, known words)
- user-specific data
- front-end-specific data
- cached data
- temporary data

It seems convenient to gather all the non-required information in one folder named *document.texd/*

document.texp such that silently removing this directory would cause no harm. As a consequence, no required data should stay inside *document.texp*, and this is the only rule concerning the required data. The *texp* file extension stands for “T_EX Project”.

In tables 4 to 9 we show the core file structure of the *document.texp* directory. This is a minimal definition involving only string encoding and spelling information because there is no consensus yet among users and all the developers of T_EX solutions, on Mac OS X at least. We make use of the XML property list data format storage as defined by <http://www.apple.com/DTDs/PropertyList-1.0.dtd>

Name	Contents
Info.plist	XML property list for any general purpose information wrapped in an info dictionary described in table 5. Optional.
<i>spellingKey.spelling</i>	XML property list for lists of known words wrapped in a spelling dictionary defined in table 9 and uniquely identified by <i>spellingKey</i> . This format is stronger than a simple comma separated list of words. Optional.
frontends	directory dedicated to front-ends only.
frontends/name	private directory dedicated to the front-end identified by <i>name</i> . The further contents definition is left under the front-end responsibility.
users	directory dedicated to users. Should not contain any front-end specific data.
users/name	directory dedicated to the user identified by <i>name</i> (not its login name). Not yet defined, but private and preferably encrypted.

Table 4: Contents of the T_EX Project directory *document.texp*

Key	Class	Contents
isa	String	Required with value: info
version	Number	Not yet used but reserved
files	Dictionary	The paths of the files involved in the project wrapped in a files dictionary. Optional.
properties	Dictionary	Attributes of the above files wrapped in a properties dictionary. Optional.
main	String	The <i>fileKey</i> of the main file, if relevant, where <i>fileKey</i> is one of the keys of the files dictionary. Optional.

Table 5: **info** dictionary description.

Key	Class	Contents
<i>fileKey</i>	String	The path of the file identified by the string <i>fileKey</i> , relative to the directory containing the T _E X project. No two different keys should correspond to the same path.

Table 6: **files** dictionary description.

However, this mechanism doesn't actually provide the concrete information needed to typeset properly (engine, format, output format). For that we can use **Makefiles** or shell scripts either embedded in the T_EX Wrapper itself or shipped as a standard tool in a T_EX distribution. This latter choice is less powerful but much more secure. Anyway, a set of default actions to be performed on a T_EX Wrapper should be outlined (compose, view, clean, archive...).

Technically, iT_EXMac uses a set of private, built-in shell scripts to typeset documents. If this is not suitable, customized ones are used instead,

but no warning is given then. No security problem has been reported yet, most certainly because such documents are not shared.

Notice that iT_EXMac declares **texd** as a document wrapper extension to Mac OS X, which means that *document.texd* folders are seen by other applications just like other single file documents, their contents is hidden at first glance. Using another file extension will prevent this Mac OS X feature without losing the benefit of the T_EX Wrapper Structure.

A final remark concerns the version control system in standard use among T_EX users. In the current definition, only one directory level should be

Key	Class	Contents
<i>fileKey</i>	Dictionary	Language, encoding, spelling information and other attributes wrapped in an attributes dictionary described in table 8. <i>fileKey</i> is one of the keys of the files dictionary.

Table 7: **properties** dictionary description.

Key	Class	Contents
isa	String	Required with value: attributes
version	Number	Not yet used but reserved
language	String	According to latest ISO 639. Optional.
codeset	String	According to ISO 3166 and the IANA Assigned Character Set Names. If absent the standard C++ locale library module is used to retrieve the codeset from the language . Optional.
eol	String	When non void and consistent, the string used as end of line marker. Optional.
spelling	String	One of the <i>spellingKeys</i> of table 4, meaning that <i>spellingKeys.spelling</i> contains the list of known words of the present file. Optional.

Table 8: **attributes** dictionary description

Key	Class	Contents
isa	String	Required with value: spelling
version	Number	Not yet used but reserved
words	Array	The array of known words

Table 9: **spelling** dictionary description.

supported in a *document.tex* folder. The contents of the **frontend** and **users** should not be monitored.

9 Nota Bene

Some features discussed here are still in the development stage and are still being tested and validated (for example, advanced syntax highlighting and full TWS support).

Šäfer \TeX : Source Code Esthetics for Automated Typesetting

Frank-Rene Schaefer

Franzstr. 21
50931 Cologne
Germany

fschaef@users.sourceforge.net
<http://safertex.sourceforge.net>

Abstract

While \TeX [4] provides high quality typesetting features, its usability suffers due to its macro-based command language. Many tools have been developed over the years simplifying and extending the \TeX interface, such as \LaTeX [5], \LaTeX3 [6], pdf \TeX [2], and NTS [8]. Front-ends such as \TeX macs [10] follow the visual/graphical approach to facilitate the coding of documents. The system introduced in this paper, however, is radical in its targetting of optimized *code appearance*.

The primary goal of Šäfer \TeX is to make the typesetting source code as close as possible to human-readable text, to which we have been accustomed over the last few centuries. Using indentation, empty lines and a few triggers allows one to express interruption, scope, listed items, etc. A minimized frame of ‘paradigms’ spans a space of possible typesetting commands. Characters such as ‘.’ and ‘\$’ do not have to be backslashed. Transitions from one type of text to another are automatically detected, with the effect that environments do not have to be bracketed explicitly.

The following paper introduces the programming language Šäfer \TeX as a user interface to the \TeX typesetting engine. It is shown how the development of a language with reduced redundancy increases the beauty of code appearance.

1 Introduction

The original role of an author in the document production process is to act as an *information source*. To optimize the flow of information, the user has to be freed from tasks such as text layout and document design. The user should be able to delegate the implementation of visual document features and styles to another entity. With this aim in mind, the traditional relationship between an author and his typesetter before the electronic age can be considered the optimal case. Modern technology has increased the speed and reduced the cost of document processing. However, the border between *information specification* and *document design* has blurred or even vanished.

In typesetting engines with a graphical user interface, an editor often takes full control over page breaks, font sizes, paragraph indentation, references and so on. Script-oriented engines such as \TeX take care of most typesetting tasks and provide high quality document design. However, quite often the task to produce a document requires detailed insight into the underlying philosophy.

Šäfer \TeX tries to get back to the basics, as depicted in Figure 1. Like the traditional writer, a user shall specify information as *redundancy-free* as possible with a minimum of commands that are alien to him. Layout, features, and styles shall be implemented according to predefined standards with a minimum of specification by the user.

To the user, the engine provides a simple interface, only requiring plain text, tables and figures. A second interface allows a human expert to adapt the engine to local requirements of style and output. Ideally, the added features in the second interface do not appear to the user, but are activated from context. Then, the user can concentrate on the core information he wants to produce, and not be distracted by secondary problems of formatting.

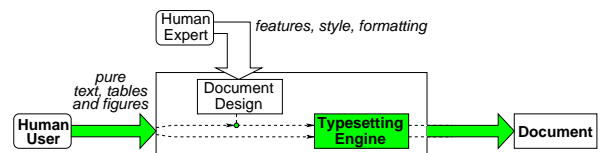


Figure 1: Neo-traditional typesetting.

The abovementioned ideal configuration of engine, user, and expert can hardly be achieved with present automated text formatting systems. While relying on TeX as a typesetting engine, ŞäferTeX tries to progress towards a minimal-redundancy programming language that is at the same time intuitive to the human writer.

2 The ŞäferTeX Engine

As shown in figure 2, the ŞäferTeX engine is based on a three-phase compilation, namely: lexical analysis, parsing and code generation. Along with the usual advantages of such modularization, this structure allows us to describe the engine in a very formal manner. In this early phase of the project, it further facilitates adding new features to the language. Using interfacing tools such as SWIG [1] and .NET [9], it would be possible in the future to pass the generated parse tree to different programming languages. Such an approach would open a whole new world to typesetters and document designers. Plugins for ŞäferTeX could then be designed in the person's favorite programming language (C++, Java, Python, C#, anything). Currently, automated document production mainly happens by preprocessing L^ATeX code. Using a parse tree, however, gives access to document contents in a structured manner, i.e., through dedicated data structures such as objects of type `section`, `item group`, and so on.

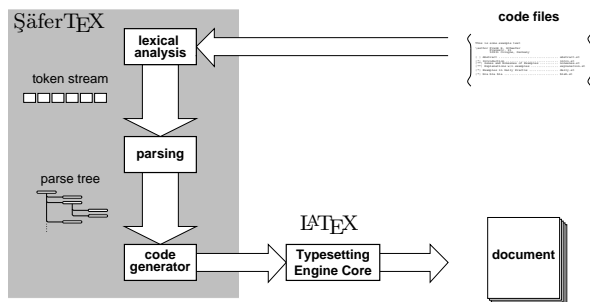


Figure 2: The ŞäferTeX compilation process.

GNU flex [7] (a free software package), is used to create the lexical analyzer. It was, though, necessary to deviate from the traditional idea of a lexical analyzer as a pure finite state automaton. A wrapper around flex implements inheritance between modes (start conditions). Additionally, the lexical analyzer produces implicit tokens and deals with indentation as a scope delimiter.

The parser is developed using the Lemon parser generator [3] (also free software). Using such a program, the ŞäferTeX language can be described with a context-free grammar. The result of the parser is

a parse tree, which is currently processed in C++. The final product is a L^ATeX file that is currently fed into the L^ATeX engine. A detailed discussion of the engine is not the intention of this paper, though. The present text focuses on the language itself.

3 Means for Beauty

ŞäferTeX tries to optimize code appearance. The author identifies three basic means by which this can be achieved:

1. The first means is intuitive treatment of characters. For example, ‘\$’ and ‘.’ are used as normal characters and do not function as commands, as they do in L^ATeX.
2. The second is to use indentation as the scope delimiter. This is reminiscent of the Python programming language. It allows the user to reduce brackets and enforces proper placement of scopes. For table environments, this principle is extended so that the column positions can be used as cell delimiters.
3. The third principle is automatic environment detection. If an item appears, then the ‘itemize’ environment is automatically assumed. This reduces redundancy, and makes the source file much more readable.

Applying these principles leads to the eight rules of ŞäferTeX as they are explained at the end (section 5). We now discuss them in more detail.

3.1 Intuitive Treatment of Characters

In the design of a typesetting language, the user has to be given the ability to enter both normal text and commands specifying document structure and non-text content. This can be achieved by defining functions, i.e., using character sequences as triggers for a specific functionality. This happens when we define, say, `\sin(x)` as a function computing the sine of `x`. For a typesetter this is not a viable option, since the character chain can be easily confused with normal text. As a result, one would have to ‘bracket’ normal text or ‘backslash’ functions. Another solution is to use extra characters. This was the method Donald Knuth chose when he designed TeX [4]. The first solution is still intuitive to most users. The second, however, is rather confusing, implying that ‘%’, ‘\$’ and ‘.’ have a meaning different from what one sees in the file.

Historically, at the time TeX was designed, keyboards had a very restricted number of characters. Moreover, ASCII being the standard text encoding in Knuth’s cultural context, the high cost of data

Table 1: Comparison of treatment of special characters in \LaTeX and $\SäferTeX$.

\LaTeX :	According to Balmun <code>\& Refish</code> <code>\$<www.b-and-r.org\$>\$</code> a conversion of module <code>\#5</code> , namely ‘ <code>propulsion_control</code> ,’ into a metric system increases code safety up to 98.7% at cost of <code>\~ \ \$17,500</code> .
$\SäferTeX$:	According to Balmun <code>&</code> <code>Refish <www.b-and-r.org></code> a conversion of module <code>#5</code> , namely ‘ <code>propulsion_control</code> ,’ into a metric system increases code safety up to 98.7% at cost of <code>~</code> <code>\$17,500</code> .

storage, and the lack of advanced programming languages also all may have contributed to the design choices made. Although the documents produced still equal and even outclass most commercial systems of our days, the input language, it must be admitted, is rather cryptic.

The first step towards readability of code is to declare a maximum number of characters as ‘normal’. In $\SäferTeX$, the only character that is not considered normal is the backslash. All other characters, such as ‘%’, ‘\$’ and ‘_’, appear in the text as they are. Special characters only act abnormal if they appear twice without whitespace in between. These tokens fall into the category of *alien things*, meaning that they look strange and thus are expected to not appear verbatim in the output.

Table 1 compares \LaTeX code to $\SäferTeX$ code, showing the improvement with respect to code appearance. The advantages may seem minor. Consider, however, the task of learning the difference between the characters that can be typed normally and others that have to be backslashed or bracketed. The abovementioned simplification already removes the chance of subtle errors appearing when \LaTeX code is compiled. The subsequent sections show how the code appearance and the ease of text input can be further improved.

3.2 Scope by Indentation

In the preceding, we discussed how commands are best defined in a typesetting engine. One way to organize information is to create specific regions, called scopes or environments. Most programming languages use explicit delimiters for scopes without giving any special meaning to white space of

Einstein clearly stated his disbelief in the boundedness of the human spirit as becomes clear through his sayings:

```
\quote The difference between genius and
      stupidity is that genius has its limits.

      Only two things are infinite, the
      universe and human stupidity, and I'm
      not sure about the former.
```

Similar reports have been heard from Frank Zappa and others.

Figure 3: Scope by indentation.

any kind. This implies that the delimiters must be visible. C++, for example, uses curly braces, while \LaTeX uses `\begin{...} ... \end{...}` constructs to determine scope. This approach allows one to place the scopes very flexibly. However, it pollutes the text with symbols not directly related to the information being described. The more scopes that are used, and the deeper they are nested, the more the source text loses readability.

Another approach is *scoping by indentation*. A scope of a certain indentation envelopes all subsequent lines and scopes as long as they have more indentation. Figure 3 shows an example of scope by indentation. \LaTeX 's redundancy-rich delimiters add nothing but visual noise to the reader of the file. $\SäferTeX$, however, uses a single backslashed command `\quote` in order to open a quote domain. The scope of the quote is then simply closed by the lesser indentation of the subsequent sentence.

This simple example was chosen to display the principle. It is easy to imagine that for more deeply nested scopes (e.g., `picture` in `minipage` in `center` in `figure`), \LaTeX code converges to unreadability, while $\SäferTeX$ code still allows one to get a quick overview about the document structure. Scope by indentation has proven to be a very convenient and elegant tool.

An extension of this concept is using *columns as cell delimiters* in a table scope. The implementation of tables in $\SäferTeX$ allows the source to omit many ‘parboxes’ and explicit ‘&’-cell delimiters. To begin with, a row is delimited by an empty line. This means that line contents are glued together as long as only one line break separates them. The cell content, though, is collected using the position of the cell markers ‘&&’ and ‘||’. Additionally, the symbol ‘~~’ glues two cells together. This makes cumbersome declarations with `\multicolumn` and

`\table Food suppliers, prices and amounts.`

Product	Price/kg	Supplier	kg	Total Price
Sugar	\$0.25	Jackie O'Neil	34	\$8.50
Yellow Swiss Cheese	\$12.2	United Independent Farmers of Switzerland	100	\$1220.00
Green Pepper Genuine Mexican	\$25.0	Anonymous Indians Tribes	2	\$50.00
Sum				\$1278.50

Figure 4: Example of writing a table: identifying cell borders by column.

`\parbox` unnecessary. Figure 4 shows an example of a ŠäferTEX table definition.

4 Implicit Environment Detection

A basic means of improving convenience of programming is reducing redundancy. In L^AT_EX, for example, the environment declarations are sometimes unnecessary. To declare a list of items, one has to specify something like

```
\begin{itemize}
  \item This is the first item and
  \item this one is the second.
\end{itemize}
```

Considering the information content, the occurrence of the `\item` should be enough to know that an itemize environment has started. Using our second paradigm, ‘scope by indentation’, the closing of the environment could be detected by the first text block that has less indentation than the item itself. The `\begin` and `\end` statements are therefore redundant. In ŠäferTEX, the token ‘--’ (two dashes) is used to mark an item. Thus, in ŠäferTEX, the item list above simply looks like:

```
-- This is the first item and
-- this one is the second.
```

As implied previously, this paradigm’s power really unfolds in combination with scope by indentation. Subsequent paragraphs simply need to be indented more than the text block to which they belong. Nested item groups are specified by higher levels of indentation, as seen in figure 5.

Some important points from the example:

- The appearance of a ‘--’ at the beginning of a line tells ŠäferTEX that there is an item and

Items provide a good means to

- structure information
- emphasize important points. There are three basic ways to do this:

[[Numbers]]: Enumerations are good when there is a sequential order in the information being presented.

[[Descriptions]]: Descriptions are suitable if keywords or key phrases are placeholders for more specific information.

[[Bullets]]: Normal items indicate that the presented set of information does not define any prioritization.

- classify basic categories

There may be other things to consider of which the author is currently unaware.

Figure 5: Example code showing ‘scope by indentation’.

that an implicit token ‘list begin’ has to be created before the token ‘item start’ is sent. The next ‘--’ signals the start of the next item.

- The ‘[[’-symbol appears at the beginning of the line. It indicates a descriptor item. Since it has a higher indentation than the ‘--’ items, it is

identified as a nested list. Therefore, an implicit token ‘list begin’ has to be created again.

- The final sentence having less indentation than anything before closes all lists, i.e., it produces implicit ‘list end’ tokens for all lists that are to be closed. Thus, the parser and code generator are able to produce environment commands corresponding to the given scopes.

The above has discussed the fundamental ideas to improve programming convenience for a typesetting system. We now turn to defining a best set of rules for expressions that implements these rules.

5 The Eight Rules of Šäfer \TeX

Rules for command design shall be consistent with the paradigms of intuitive treatment of characters, scope by indentation, and automatic environment detection. The following set of rules was designed to meet these goals for Šäfer \TeX while striving for a intuitive code appearance:

[1] Every character and every symbol in the code appears in the final output as in the source document, except for ALIEN THINGS.

[2] ALIEN THINGS look alien.

In plain \TeX , characters such as ‘\$’, ‘%’ and ‘_’ do not appear in the document as typed. The fact that they look natural but trigger some \TeX specific behavior is prone to confuse the layman. In Šäfer \TeX , they appear as typed on the screen. Alien things can be identified by their look. The next four rules define the ‘alien look:’

[3] Any word starting with a single backslash \backslash . Examples are `\figure` and `\table`.

[4] Any non-letter character that appears twice or more, such as ‘##’ (this triggers the start of an enumeration item at the beginning of the line).

[5] Parentheses (at the beginning of a line) that only contain asterisks ‘*’ or whitespace. Sequences such as ‘(*)’, ‘()’, ‘(***)’ indicate sections and subsections.

[6] The very first paragraph of the file. It is interpreted as the title of the document.

Except for the first case, alien things do not interfere with readability. In fact, the double minus ‘--’ for items and the ‘(*)’ for sections are used naturally in many ASCII files. Internally, alien things are translated into commands for the typesetting engine, but the user does not need to know.

The last two issues are separation of the text stream and identification of scope of an environment:

[7] Termination of paragraphs, interruptions of the text flow, etc., are indicated by an EMPTY LINE.

[8] The scope of an environment, table cells, etc. is determined by its INDENTATION. A line with less indentation closes all scopes of higher indentation.

These are the eight rules of Šäfer \TeX which enable one to operate the typesetter. They are defined as ‘rules’ but, in fact, they do not go much beyond common organization of text files.

6 Commands

This section gives a brief overview of the commands that are currently implemented. In this early stage of development, the system’s structure and language design has been in the foreground, in order to build the framework for a more powerful typesetting engine. In the current version of Šäfer \TeX , the following commands are implemented:

--, ++ starts a bullet item. The two can be used interchangeably to distinguish different levels of nested item groups.

starts an enumeration item.

[[]] bracket the beginning of a description item.

`\table` opens a table environment. It is followed by a caption and the table body as described in section 3.2.

`\figure` opens a figure environment. The text following this command is interpreted as the caption. Then file names of images are to be listed. Images that are to be shown side by side are separated by ‘&&’. Vertically adjacent images are separated by empty lines.

`\quote` opens a quote environment.

(*) starts a section. The number of asterisks indicates the level of the section.

() starts a section without a section number. The number of blanks indicates the section level.

.... includes a file (more than four dots is equivalent to four). The next non-whitespace character sequence is taken as the filename to be included.

`\author` specifies information about the author of the document.

Commands have been designed for footnotes, labels, and more. However, due to the early stage of development, no definite decision about their format has been made. In the appendix, two example files are listed in order to provide an example of Šäfer \TeX code in practical applications.

7 Conclusion and Outlook

Using simple paradigms for improving code appearance and reducing redundancy, a language has been developed that allows more user-friendly input than is currently possible with T_EX and L^AT_EX. These paradigms are: the intuitive processing of special characters, the usage of indentation for scope and the implicit identification of environments. As an implementation of these paradigms, the eight rules of ŠäferT_EX were formed, which describe the fundamental structure of the language.

While developing ŠäferT_EX, the author quickly realized that the ability to provide the parse tree to layout designers extends the usage beyond the domain of T_EX. Currently, much effort remains to provide appropriate commands for document production. Functionality of popular tools such as `psfrag`, `fancyheaders`, `bibtex`, `makeindex`, etc., are to be implemented as part of the language. In the long run, however, it may be interesting to extend its usage towards a general markup language.

8 Acknowledgments

The author would like to thank the T_EX Stammtisch of Cologne in Germany for their valuable comments. Special thanks to Holger Jakobs who helped translate this text from Genglish to English.

References

- [1] D. M. Beazley. Automated scientific software scripting with SWIG. In *Tools for program development and analysis*, volume 19, pages 599–609. Elsevier Science Publishers B. V., Amsterdam, The Netherlands, 2003.
- [2] T. T. H^àn, S. Rahtz, and H. Hagen. The pdftex manual. <http://www.ntg.nl/doc/han/pdftex-a.pdf>, 1999.
- [3] R. D. Hipp. The Lemon Parser Generator. <http://www.hwaci.com/sw/lemon>, 1998.
- [4] D. E. Knuth. *The T_EXbook*. Addison Wesley, 1983.
- [5] H. Kopka and P. Daly. *A Guide to L^AT_EX*. Addison Wesley, 1992.
- [6] Frank Mittelbach and Chris Rowley. The L^AT_EX3 Project. *TUGboat*, 18(3):195–198, 1997.
- [7] J. Poskanzer and V. Paxson. Flex, a fast lexical analyzer generator. <http://sourceforge.net/projects/lex>, 1995.
- [8] P. Taylor, J. Zlatuška, and K. Skoupy. *The NTS project: from conception to implementation*. Cahiers GUTenberg, May 2000.
- [9] A. Troelsen. *C# and the .NET Platform*. APress, 2001.
- [10] Joris van der Hoeven. GNU T_EXmacs. <http://www.texmacs.org/tmweb/home/welcome.en.html>, 2003.

Details about
The Elves and The Shoemaker

```

\Author Original:
  Brothers Jakob & Wilhelm Grimm
  Somewhere in Germany

( ) Abstract ..... abstract.st

(*) Nocturne shoe productions ..... strange.st
(**) Living in confusion ..... confusion.st
(**) Women make trouble ..... trouble.st

(*) Midnight observations ..... midnight.st
(**) Elves in the cold ..... freezing-elves.st

(*) New era for elves: luxury ..... luxury.st

(*) Elves leave their job undone ..... spoiled-elves.st
  
```

Figure 6: Example input ‘main.st’.

```

\figure ::fig:plots:: Performance a) productivity of shoemaker. b) gain.

  ferry-tales/prod.eps  &&  ferry-tales/capital.eps

Reviewing the plots of shoes produced (figure --<fig:plots>), the shoemaker
realized an instantaneous increase during the night period. He only could
think of two possible reasons:

## He was sleepworking. Since he even used to work few when awake this
  assumption was quickly refuted.

## Elves must have come over night and did some charity work.

He further based his theory on the influence of the tanning material used. In
fact, there were differences in the number of shoes produced depending on acid
number and pH value (see table --<tab:tan-mat>).

\table ::tab:tan-mat::Influence of tanning materials on shoe production.

  Tanning Mat. && pH value  && acid number && shoes prod.
-----
  European   && 3.4 - 3.7 && 30 - 40 && 32      @@
  Indian     && 2.0 - 2.1 && 31 - 45 && 35      @@
  African    && 4.5 - 4.6 && 33 - 37 && 36      @@
  Australian && 3.0 - 7.0 && 27 - 45 && 15      @@
-----
Resourcing several leathers from indian & african suppliers allowed him to
increase profit ranges tremendously. Moreover, these shoes were sold at an
even higher price around $0.50. Pretty soon, the shoemaker was able to save a
good sum of $201.24.
  
```

Figure 7: Example input ‘confusion.st’.

FEATPOST and a Review of 3D METAPOST Packages

L. N. Gonçalves

CFMC-UL, Av. Prof. Gama Pinto 2

1649-003 Lisboa

Portugal

nobre@lince.cii.fc.ul.pt

<http://matagalatlante.org>

Abstract

METAPOST is able to produce figures that look almost like ray-traced raster images but that remain vector-based. A small review of three-dimensional perspective implementations with METAPOST is presented. Special emphasis is given to the abilities of the author's implementation: FEATPOST.

1 Introduction

There are at least four METAPOST packages related to three-dimensional diagrams:

- GNU 3DLDF — <http://www.gnu.org/directory/graphics/3D/3DLDF.html>
- 3d/3dgeom — <http://tug.org/tex-archive/graphics/metapost/macros/3d/>
- m3D — <http://www-math.univ-poitiers.fr/~phan/m3Dplain.html>
- FEATPOST — <http://matagalatlante.org/nobre/featpost/doc/featexamples.html>

All of these packages are individual and independent works “under construction”. There has been neither collaboration nor competition among the authors. Each produces different kinds of diagrams and each uses a different graphic pipeline. The following sections of this document describe these packages, in a mainly independent way.

2 GNU 3DLDF

3DLDF is not a pure METAPOST package, as it is written in C++ using CWEB. Diagrams are also coded in C++ and are compiled together with the package. Nevertheless, this is, of all four, the package with the greatest promise for a future three-dimensional-capable METAPOST.

1. It outputs METAPOST.
2. Its syntax is similar to METAPOST.
3. It overcomes the arithmetic limitations inherent in METAPOST.
4. Both the affine transformations and the graphics pipeline are implemented through 4×4 matrices.
5. Its author, Laurence D. Finston, is actively improving and maintaining the package. His plan

includes, among many other ideas, the development of an input routine (to allow interactive use) and the implementation of three-dimensional paths via NURBS.

Given the possible computational efficiency of this approach, one can foresee a system that merges the METAPOST language with the capabilities of standard ray-tracing software.

3 3d/3dgeom

This was the first documented extension of METAPOST into the third dimension—and also into the fourth dimension (time). Denis B. Roegel created, back in 1997, the 3d package to produce animations of polyhedra. In 2003 he added the 3dgeom “module” which is focused on space geometry. It remains the least computationally intensive package of those presented here.

1. Each component of a point or a vector is stored in a different numeric array. This eases control of a stack of points. Points are used to define planar polygons (faces of polyhedra) and the polygons are used to define *convex* polyhedra.
2. When defining a polygon, a sequence of points must be provided such that advancing on the sequence is the same as rotating clockwise on the polygon, when the polygon is visible. This means that, when a polyhedron is to be drawn, the selection of polygons to be drawn is very easy: only those whose points rotate clockwise (the visible ones). Hidden line removal is thus achieved without sorting the polygons.
3. Points can also be used to define other points according to rules that are common in the geometry of polyhedra or according to operations involving straight lines and/or planes and/or angles.

4. The author plans to release an updated version with the ability to graph parametric lines and surfaces.

4 m3D

Anthony Phan developed this very interesting package but has not yet written its documentation. Certainly, this is, of all four, the package that can produce the most complex and beautiful diagrams. It achieves this using, almost exclusively, four-sided polygons.

1. Complex objects can be defined and composed (see figure 1). For example, one of its many predefined objects is the fractal known as the “Menger Sponge”.
2. It can render revolution surfaces defined from a standard METAPOST path (see figure 2).
3. Objects or groups of polygons can be sorted and drawn as if reflecting light from a punctual source and/or disappearing in a foggy environment.

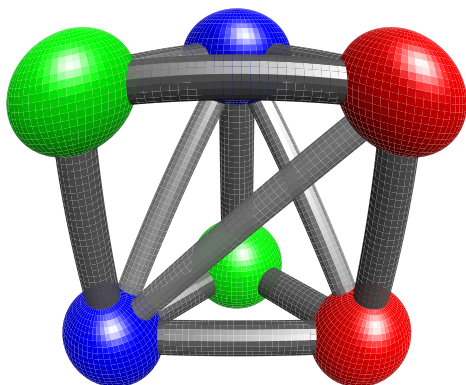


Figure 1: A diagram produced by m3D showing a single object, composed of spheres and cylindrical connections, under a spherical perspective.

5 FEATPOST

Geared towards the production of physics diagrams, FEATPOST sacrifices programming style and computational efficiency for a large feature set.

1. Besides the usual parallel and central perspectives it can make a sort of “spherical distortion” as if a diagram is observed through a fish-eye lens¹. This kind of perspective is advantageous for animations as it allows the point of view to be inside or among the diagram objects. When using the central perspective, points that are as

¹ Also possible with m3D.

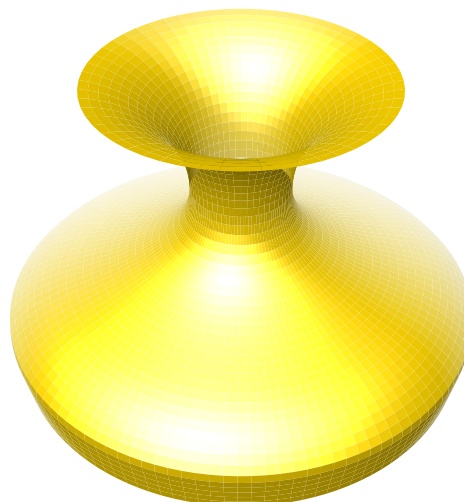


Figure 2: A diagram produced by m3D showing a revolution surface under a central perspective.

distant from the projection plane as the point of view get projected at infinity, and MetaPost overflows and crashes. The spherical projection is always finite.

2. It can mark and measure angles in space.
3. It can produce shadows of some objects (see figure 9). Shadows are calculated in much the same way as perspectives. The perspective projection, from 3D into 2D, is a calculation of the intersection of a straight line and a plane. A shadow is also a projection from 3D into 2D, only the line and the plane are different. The shadow must be projected onto the paper page before the object that creates the shadow. Shadows are drawn after two projections, objects are drawn after one projection and after their shadows.
4. It can correctly draw intersecting polygons (see figure 12).
5. It knows how to perform hidden line removal on some curved surface objects. Imagine a solid cylinder. Now consider the part of the cylinder’s base that is the farthest away. You only see a part of its edge. In order to draw that part, it is necessary to know the two points at which the edge becomes hidden. FEATPOST calculates this. Note that the edge is a circle, a curved line. FEATPOST does not use polygons to hide lines on some curved surface objects.
6. Supported objects include: dots, vectors, angles, ropes, circles, ellipses, cones, cylinders, globes, other curved surface objects, polygons, cuboids, polyhedra, functional and parametric

surface plots, direction fields, field lines and trajectories in conservative force fields.

Many of the drawable objects are not made of polygons, but rather of two-dimensional paths. FEATPOST does not attempt to draw surfaces of these objects, only their edges. This is partly because of the use of intrinsic METAPOST functions and partly because it eases the production of diagrams that combine space and planar (on paper) objects.

One of the intrinsic METAPOST functions that became fundamental for FEATPOST is the composition `makepath makepen`. As this converts a path into its convex form, it very much simplifies the determination of some edges.

Another important aspect of the problem is hidden line removal. Hidden line removal of a group of polygons can, in some cases, be performed by drawing the polygons by decreasing order of distance to the point of view. FEATPOST generally uses the Shell sorting method, although when the polygons are just the faces of one cuboid FEATPOST has a small specific trick. There is also a specific method for hidden line removal on cylinders and another for other curved surface objects.

5.1 Examples

Some of the FEATPOST macros are presented here. Detailed information is available at

- <http://matagalatlante.org/nobre/featpost/doc/macroMan.html>
- <http://www.ctan.org/tex-archive/graphics/metapost/macros/featpost/>

Each perspective depends on the point of view. FEATPOST uses the global variable `f`, of type `color`, to store the (X, Y, Z) coordinates of the point of view. Also important is the aim of view (global variable `viewcentr`). Both together define the line of view.

The perspective consists of a projection from space coordinates into planar (u, v) coordinates on the projection plane. FEATPOST uses a projection plane that is perpendicular to the line of view and contains the `viewcentr`. Furthermore, one of the projection plane axes is horizontal and the other is on the intersection of a vertical plane with the projection plane. “Horizontal” means parallel to the XY plane.

One consequence of this setup is that `f` and `viewcentr` must not be on the same vertical line (as long as the author avoids solving this problem, at least!). The three kinds of projection known to FEATPOST are schematized in figures 3, 4 and 5.

The macro that actually does the projection is, in all cases, `rp`.

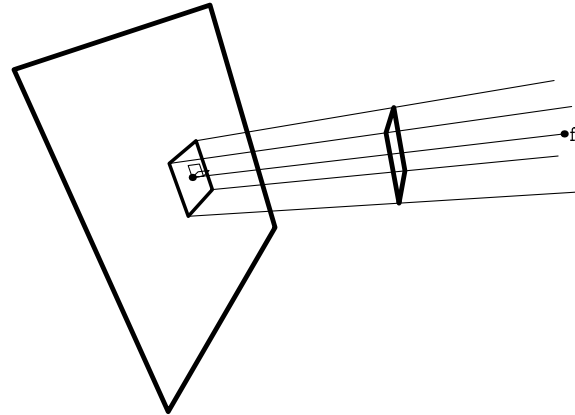


Figure 3: Parallel projection.

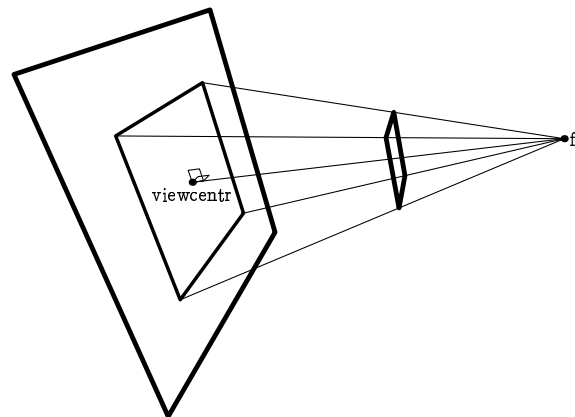


Figure 4: Central projection.

Physics problems often require defining angles, and diagrams are needed to visualize their meanings. The `angline` and `squareangline` macros (see figure 6 and the code below) support this.

```
f := (5,3.5,1);
beginfig(2);
  cartaxes(1,1,1);
  color va, vb, vc, vd;
  va = (0.29,0.7,1.0);
  vb = (X(va),Y(va),0);
  vc = N((-Y(va),X(va),0));
  vd = (0,Y(vc),0);
  drawarrow rp(black)--rp(va);
  draw rp(black)--rp(vb)--
      rp(va) dashed evenly;
  draw rp(vc)--rp(vd) dashed evenly;
  drawarrow rp(black)--rp(vc);
  squareangline( va, vc, black, 0.15 );
```

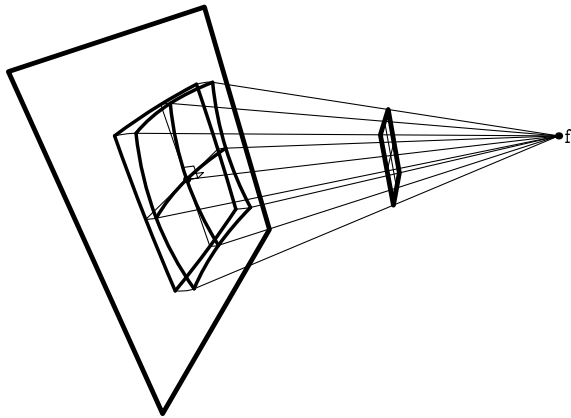


Figure 5: Spherical projection. The spherical projection is the composition of two operations: (i) there is a projection onto a sphere and (ii) the sphere is plaited onto the projection plane.

```

angline(va,red,black,0.75,
        decimal getangle(va,red),lft);
endfig;

```

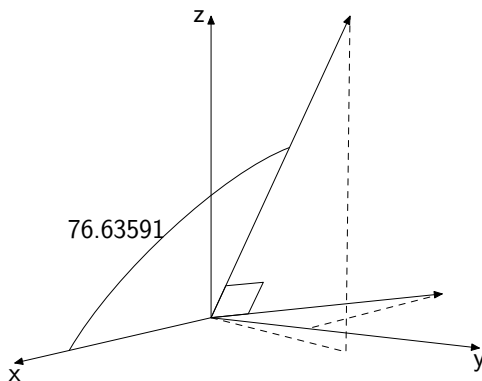


Figure 6: FEATPOST diagram using `angline`.

Visualizing parametric lines is another need of physicists. When two lines cross, one should be able to see which line is in front of the other. The macro `emptyline` can help here (see figure 7 and the code below).

```

f := (2,4,1.8);
def theline( expr TheVal ) =
  begingroup
    numeric cred, cgre, cblu, param;
    param = TheVal*(6*360);
    cred = -0.3*cosd( param );
    cblu = 0.3*sind( param );
    cgre = param/850;
    ( cred,cgre,cblu )
  endgroup

```

```

enddef;
beginfig(1);
  numeric axsize, zaxpos, zaxlen;
  color xbeg, xend, ybeg,
        yend, zbeg, zend;

  axsize = 0.85;
  zaxpos = 0.55;
  zaxlen = 2.1;
  pickup pencircle scaled 1.5pt;
  xbeg = (axsize,0,0);
  xend = (-axsize,0,0);
  ybeg = (0,0,-axsize);
  yend = (0,0,axsize);
  zbeg = (zaxpos,-zaxpos,0);
  zend = (zaxpos,zaxlen,0);
  drawarrow rp( xbeg )--rp( xend );
  drawarrow rp( ybeg )--rp( yend );
  defaultscale := 1.95;
  label.rtl( "A", rp( xend ) );
  label.lft( "B", rp( yend ) );
  emptyline(false,1,black,
            0.5black,1000,0.82,2,theline);
  drawarrow rp( zbeg )--rp( zend );
  label.bot( "C", rp( zend ) );
endfig;

```

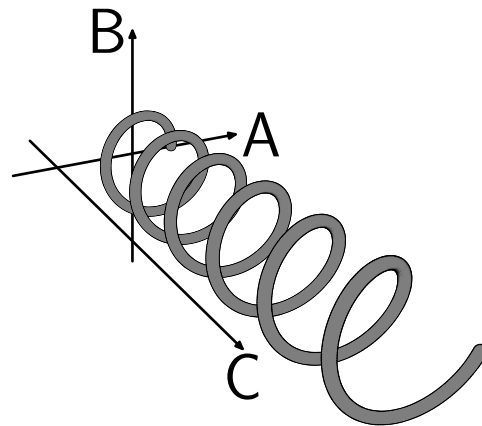


Figure 7: FEATPOST diagram using `emptyline`.

Cuboids and labels are always needed. The macros `kindofcube` and `labelinspace` fulfill this need (see figure 8 and the code below). The macro `labelinspace` does not project labels from 3D into 2D. It only Transforms the label in the same way as its bounding box, that is, the same way as two perpendicular sides of its bounding box. This is only exact for parallel perspectives.

```

f := (2,1,0.5);
ParallelProj := true;
verbatimtex

```



```

\documentclass{article}
\usepackage{beton,concmath,ccfonts}
\begin{document}
etex
beginfig(1);
  kindofcube(false,true,(0,-0.5,0),
             90,0,0,1.2,0.1,0.4);
  kindofcube(false,true,(0,0,0),
             0,0,0,0.5,0.1,0.8);
  labelinspace(false,(0.45,0.1,0.65),
              (-0.4,0,0),(0,0,0.1),
              btex
              \framebox{\textsc{Label}}
              etex);
endfig;
verbatim \end{document} etex

```

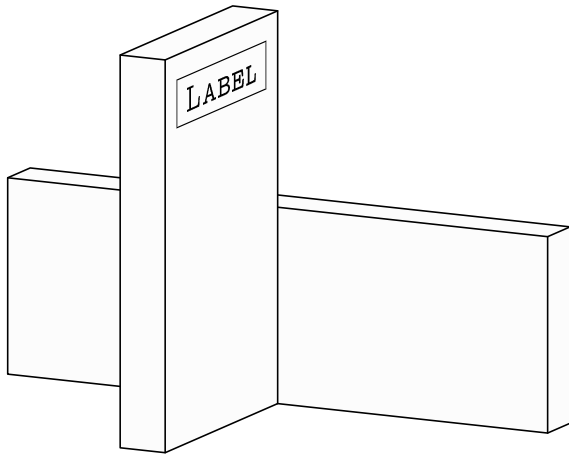


Figure 8: FEATPOST diagram using the macros `kindofcube` and `labelinspace`.

Some curved surface solid objects can be drawn with FEATPOST. Among them are cones (`verygoodcone`), cylinders (`rigorousdisc`) and globes (`tropicalglobe`). These can also cast their shadows on a horizontal plane (see figure 9 and the code below). The production of shadows involves the global variables `LightSource`, `ShadowOn` and `HoriZon`.

```

f := (13,6,4.5); ShadowOn := true;
LightSource := 10*(4,-3,6);
beginfig(3);
  numeric refln, frac, coordg;
  numeric fws, NumLines;
  path ella, ellb;
  color axe, cubevertex, conecenter,
        conevertex, allellaxe, ellaaxe,
        pca, pcb;
  frac := 0.5;          wang := 60;

```

```

axe := (0,cosd(90-wang),
        sind(90-wang));
fws := 4;          refln := 0.35*fws;
coordg := frac*fws;
NumLines := 45;
HoriZon := -0.5*fws;
setthestage(0.5*NumLines,3.3*fws);
cubevertex = (0.3*fws,-0.5*fws,0);
tropicalglobe( 7, cubevertex,
              0.5*fws, axe );
allellaxe:=reflen*(0.707,0.707,0);
ellaaxe:= refln*( 0.5, -0.5, 0 );
pcb := ( -coordg, coordg, 0 );
rigorousdisc( 0, true, pcb,
              0.5*fws, -ellaaxe );
conecenter =
  ( coordg, coordg, -0.5*fws );
conevertex = conecenter +
  ( 0, 0, 0.9*fws );
verygoodcone(false,conecenter,
             blue,reflen,conevertex);
endfig;

```

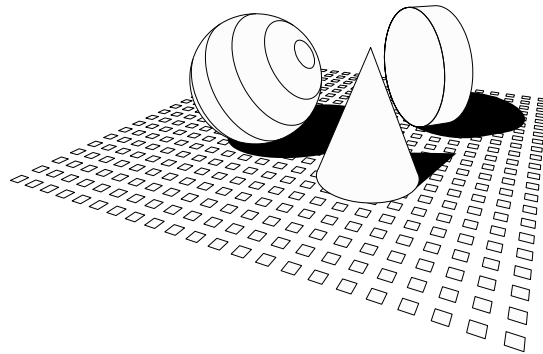


Figure 9: FEATPOST diagram using the macros `rigorousdisc`, `verygoodcone`, `tropicalglobe` and `setthestage`.

Another very common need is the plotting of functions, usually satisfied by software such as Gnuplot (<http://www.gnuplot.info/>). Nevertheless, there are always new plots to draw. One kind of FEATPOST plot that just became possible is the “triangular grid triangular domain surface” (see figure 10 and this code):

```

f := 16*(4,1,1);
LightSource := 10*(4,-3,4);
def zsu( expr xc, yc ) =
  cosd(xc*57)*cosd(yc*57)+
  4*mexp(-(xc**2+yc**2)*6.4) enddef;
beginfig(1);
  hexagonaltrimesh(false,52,15,zsu);
endfig;

```

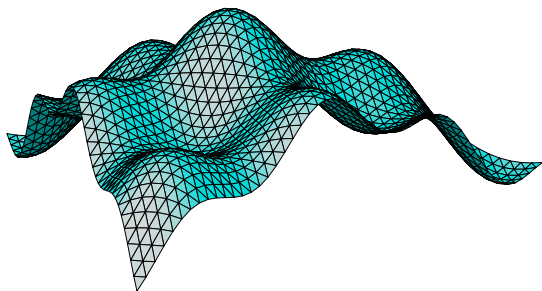


Figure 10: FEATPOST surface plot using the macro `hexagonaltrimesh`.

One feature that merges 2D and 3D involves what might be called “fat sticks”. A fat stick resembles the Teflon magnets used to mix chemicals. They have volume but can be drawn like a small straight line segment stroked with a `pencircle`. Fat sticks may be used to represent direction fields (unitary vector fields without arrows). See figure 11 (the code is skipped from now on).

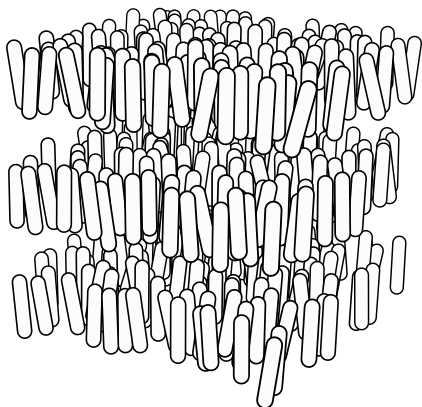


Figure 11: FEATPOST direction field macro `director_invisible` was used to produce this representation of the molecular structure of a Smectic A liquid crystal.

Finally, it is important to remember that some capabilities of FEATPOST, although usable, may be considered “buggy” or only partially implemented. These include the calculation of intersections among polygons, as in figure 12, and the drawing of toruses, as in figure 13. These two figures show “usable” situations but their code is skipped.

FEATPOST has many macros: some are specifically for physics diagrams, others may be useful for general purposes, some do not fit in this article and,

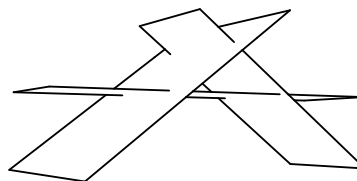


Figure 12: Intersecting polygons drawn with the macro `sharpraytrace`.

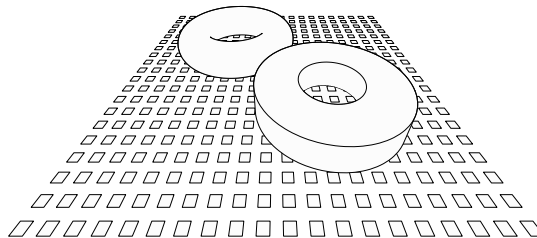


Figure 13: Final FEATPOST example containing a `smoothertorus` and a `rigorousdisc` with a hole. These macros may fail for some view points.

sadly, some are not anywhere documented. For instance, the tools for producing animations are not yet documented. (These tools are completely external to \TeX : the control of an animation is done with a Python script, and Ghostscript and `netpbm` are used to produce MPEG videos.)

In summary, the collection of three-dimensional METAPOST software, such as the four reviewed packages, is large and growing in many independent directions. It constitutes an excellent resource for those desiring to produce good diagrams.

6 Acknowledgements

Many people have contributed to make FEATPOST what it is today. Perhaps it would have never come into being without the early intervention of Jorge Bárrios, providing access to his father’s computer. Another fundamental moment happened when José Esteves first spoke about METAPOST.

More recently, the very accurate criticism of Cristian Barbarosie has significantly contributed to the improvement of these macros. Jens Schwaiger contributed new macros. Pedro Sebastião, João Dinis and Gonçalo Morais proposed challenging new features. The authors of the other packages graciously reviewed the paper, and Karl Berry actually entered new text into this document. They all have my deep thanks.

The bigfoot Bundle for Critical Editions*

David Kastrup
Kriemhildstr. 15
44793 Bochum
Germany
dak@gnu.org

Abstract

The L^AT_EX package `bigfoot` and supporting packages solve many of today's problems occurring in the contexts of single and multiple blocks of footnotes, and more. The main application is with philological works and publications, but simpler problems can be solved painlessly as well without exercising all of the package's complexities. For other problems not yet tackled in this area, a solid framework is provided.

1 Introduction

Last year, the author was approached about creating the necessary L^AT_EX style files for typesetting a critical edition of the complete works of Ernst Troeltsch¹.

With the typical optimism² that is customary among programmers, the task was accepted. "Thus, `bigfoot` was born" would be an exaggeration since it only came into being after quite a few attempts failed.³ The main reason for failure was the idea that one might preprocess nested insertions in a way that would make T_EX's own insertion splitting routines⁴ serve a useful purpose.

So let us concentrate on the present survivor instead. Some of its features are due to the original requirements, some of them are due to the author wanting to provide them in spite of not being needed by the current project. I am grateful to TUG for sponsoring some of those aspects that turn the package into something more generally useful for the T_EX community.

2 Features

So what are the features that `bigfoot` provides?

- Multiple footnote apparatus⁵ are possible.⁶

* and a lot of other footnote applications

¹ a theologian of the last century

² read: mis-estimate of work and time ³ mostly because of unmanageable^a complexity ⁴ expletive deleted

⁵ An apparatus is one block of contiguous footnotes forming a logical and physical unit. Separate apparatus^b can be independently broken to the next page.

⁶ Actually, `manyfoot` already provides this functionality^c

^a Well, for me ^b Yes, this is the correct plural form.

^c and is loaded by `bigfoot`

- Footnotes can be nested.⁷
- Footnotes are numbered in the order they appear on the page, and numbering may start from 1[†] on each page. In each apparatus, the footnotes are arranged in numerical order identical to page order. This does not sound exciting at all until you consider the implications of footnotes being nested: if the main text has some footnote⁸ and then the publisher comments the main text with a footnote,^d the logical order of footnotes (in which they appear in the source text) would have been to let footnote f appear before footnote d. The footnotes instead will be reordered to page order.⁹
- Footnotes may contain `\verbatim` commands¹⁰ and similar, and they will just work as expected. This is achieved in a manner similar to the `\footnote` command of plain T_EX.
- Footnotes can be broken across pages.¹¹

but it fails to address a number of intricacies inherent to this sort of setup, a few of which follow.

⁷ You can anchor footnotes for some apparatus in the main text^e.

[†] or whatever the first footnote symbol happened to be

⁸ such as shown in this example footnote^f

⁹ The style file `perpage` has been extended with additional functionality for reordering such numbers.

¹⁰ even stuff like `\verb-\iffalse-`

¹¹ While this does not sound like something excitingly new, it must be noted that T_EX does not do a satisfactory job at splitting insertions, the underlying mechanism for split footnotes. In particular, T_EX only manages to find a split when

^d This is a subsequent comment to the main text.

^e or any apparatus preceding it on the page

^f which happens to have a comment attached to it. Notice that `bigfoot` will prefer to leave this smaller footnote block intact, as breaking it will not help fitting the above footnote block on the page.

- When footnotes are broken across pages, the color stack is maintained properly. Color is handled in \LaTeX with the help of specials that switch the color (and, in the case of `dvips`, restore it afterwards with the help of a color stack). Restarting the footnote on the next page with the proper color is something that has never worked in \LaTeX . Now it simply does.
- Footnotes may be set in a compact form in one running paragraph.¹²

no material whatsoever is added to the page after the occurrence of the split footnote. This might include another footnote in a different apparatus, or simply a line tied to the current line with an infinite penalty, for example because of a respective setting of `\widowpenalty`. In contrast, `bigfoot` breaks footnotes properly in such circumstances, and it uses a backtracking algorithm (with early pruning of branches that can't beat the current optimum) for finding the best split positions for several footnote apparatus in parallel. The fill level of the page is taken into account as well as the costs of the individual splits. A split footnote is penalized with a penalty of 10000 (which is pretty similar to what \TeX itself does when dealing with footnotes), so that in general \TeX will tend to avoid splitting more than a single footnote whenever possible. One complication is that if the parts broken to the next page contain footnotes themselves, those have to be moved to the next page completely and adapted to the numbering of footnotes there^a. This rather intricate and complicated mechanism leads to results that look simple and natural.

¹² While `manyfoot` and `fnpara` also offer this arrangement, `bigfoot` offers a superior solution in several respects:

- The line breaking can be chosen much more flexibly: with appropriate customization, it is possible to fine-tune quite well when and where stuff will be placed in the same line, and when starting a new line will be preferred.
- In-paragraph footnotes can be broken across pages automatically, just like normal footnotes. They will only be broken after the last footnote in the block has started.
- Pages will not become over- or underfull because of misestimating the size of in-paragraph footnotes. Also the total width of such footnotes is not restricted to `\maxdimen` (which sounds generous at something like 6 m or 19 ft, until you realize that a few pages of text suffice to burst that limit, and a few pages of text are reached easily with longer variants of the main text). While \TeX will accumulate boxes exceeding this size without problem, it panics at its own audacity if you actually ask about the total width of the acquired material. While one may still not have material exceeding a total *vertical* size of `\maxdimen` accumulate in one footnote block, one would usually need a few dozen pages for that, and so *this* limitation is much less noisome than the corresponding restriction on the horizontal size.
- The decision of whether a footnote in-paragraph or standalone can be changed for each footnote apparatus at any time, including on mid-page. In fact, you can make this decision for each footnote separately. Since display math requires vertical mode footnotes, this is convenient.

^a which can be completely different!

- Split footnotes will not get jumbled in the presence of floats. `bigfoot` is not afflicted by this bug in \LaTeX 's output routine since it does not delegate the task of splitting footnotes to \TeX in the first place. While the faulty output routine of \LaTeX may still jumble the order of footnotes in that particular case (when one footnote gets held over as an insertion 'floated' at infinite cost), `bigfoot` will sort the jumbled footnotes back into order before processing them.
- Each footnote apparatus can have its own private variant of `\@makefntext` and a few other macros and parameters responsible for formatting a footnote block. The default is to use what the class provides, but special versions can be defined, for example,

```
\FootnoteSpecific{variants}%
\long\def\@makefntext#1{...
```

for the footnote block called "variants".

3 Drawbacks

What about current drawbacks?

- ε - \TeX is used throughout. After it became clear that the implementation of the package would not be possible without using some of ε - \TeX 's features, its features were extensively employed: rewriting the package to get along without ε - \TeX would be very hard, even if you came up with ideas for those cases where I could find no other solution. Free \TeX distributions have come with ε - \TeX for a long time by now (in fact, ε - \TeX is now the recommended engine for \LaTeX , and actually used as the default in the latest \TeX Live), but proprietary variants may lack ε - \TeX support. The same holds for quite a few Ω versions.
- The licence is not the LPPL, but the GPL. In my book, I consider this an advantage: the functionality of the package is quite important, and it is in its infancy yet. I would not like to encourage a market of proprietary offspring directly competing with it. While with sufficient financial incentive I might feel confident enough to have the means to reimplement whatever noteworthy extension somebody else might come up with, at the current time I prefer this way of ensuring that the free development does not fall behind and that there is no incentive to

-
- `bigfoot` will make a good-faith effort to adapt the normal footnote layout provided by the document class with the `\@makefnmark` and `\@makefntext` macros to in-paragraph footnotes.

turn to developers with no qualms about creating proprietary versions.

- **bigfoot** requires twice as many box registers¹³ as **manyfoot**: one set in the form of an insertion for each footnote apparatus, one set as mere boxes.
- It can't handle more footnotes in a single block per page than the group nesting limit of **TeX**, and that is usually hardwired at 255.[†]
- Since it meddles considerably with the output routine's workings, interoperation with other packages doing the same might be problematic. Considerable effort has been spent on minimizing possibly bad interactions, but the results might not always be satisfactory and, at the very least, might depend on the load order of packages.
- It slows things down. This is not much of a concern, and usually the package is astonishingly fast.
- The complexity of the package makes it more likely for things to go wrong in new ways.¹⁴

4 Additional New Packages

The bundle provides some more packages: **perpage** is used for the sort of renumbering games mentioned before, and **suffix** is used for defining augmented commands.

As an example of use for those packages we had previously a few examples where numbers like 7[‡] and 255[§] were given footnotes, and in order not to confuse this with powers as the following 666¹⁵ is in danger of, we have switched to per-page numbering of footnotes with symbols for that purpose. The source code simply uses

```
like~7\footnote' {a lucky number}
```

namely a variant footnote command. How is that achieved? Just with

¹³ Since ϵ -**TeX** has an ample supply of box registers (32767 instead of 256), this is not really much of an additional limitation. If you find yourself running out of insertions, **etex.sty** offers the `\reserveinserts` command.

[†] This limit seems sufficient at first glance, but one could use the various mechanisms available in connection with in-paragraph footnotes to make sure that a footnote will be broken across the page at a point closely related to the main text's breakpoint (for example, if you are doing an interlinear translation in a footnote). In that case, this limit might become problematic.

¹⁴ Most of those problems should arise under requirements that could not possibly be met without the package, so this would be reason for improving rather than not using the package.

[‡] a lucky number [§] well, almost as lucky

¹⁵ strange, yes?

```
\newcounter{footalt}
\def\thefootalt{\fnsymbol{footalt}}
\MakeSortedPerPage[2]{footalt}
\WithSuffix\def\footnotedefault' {%
  \refstepcounter{footalt}%
  \Footnote{\thefootalt}}
```

A new counter is created, its printed representation is set to footnote symbols, the counter is made to start from 2 on each page (since symbol 1[¶] is a bit ugly), and then a variant of `\footnotedefault` is defined which will step the given counter and use it as a footnote mark.¹⁶

That's all. One can define several suffixes, the resulting commands are robust¹⁷, and one can use arguments and other stuff. For example,

```
\WithSuffix\long\def\footnotedefault
[#1]{#2}{...}
```

would augment the macro `\footnotedefault` by a variant accepting an optional argument.

5 Some Internals

5.1 Basic Operation

The package uses most of the interfaces of **manyfoot** for its operation. While it uses **TeX**'s insertions for managing the page content, the material collected in those insertions is in a pretty raw state and its size is always overestimated.¹⁸ The actual material that goes onto the finished page is generated from the insertions at `\output` time.

Material that is put into insertions is prewrapped into boxes without intervening glue.¹⁹ The box dimensions are also somewhat special: while the total height (height+depth) corresponds to the actual size of the footnote, the depth contains a unique id that identifies the last footnote in each box (of which there usually is just one, unless we are dealing with the remnants of an in-paragraph footnote apparatus broken across pages). The width is set to a sort key that is used for rearranging the various footnotes into an order corresponding to their order of appearance on the page.

[¶] which is *

¹⁶ **manyfoot** defines a two-argument command `\Footnote` that takes a footnote mark and corresponding footnote text.

¹⁷ as long as their suffixes are so as well

¹⁸ **bigfoot** simply sets each footnote, even those that should be typeset with others in one block, separately in its own paragraph for estimating its size, which should be a safe upper limit for the size a footnote can take when set in a paragraph with others.

¹⁹ That way, there is never a legal breakpoint in an insertion.

The boxes are sorted by unboxing them and then calling the comparatively simple sorting routine (a straight insertion sort):

```
\def\FN@sortlist{%
  \setbox\z@\lastbox
  \ifvoid\z@ \else
    \FN@sortlist \FN@sortlistii
  \fi}}
```

```
\def\FN@sortlistii{%
  \setbox\tw@\lastbox
  \ifvoid\tw@\else
    \ifdim\wd\tw@<\wd\z@
      {\FN@sortlistii}%
    \fi
    \nointerlineskip \box\tw@
  \fi
  \nointerlineskip \box\z@}
```

and then all consecutive runs of hboxes are joined into vboxes. The desirability of breaking between two in-paragraph footnotes depends on their respective size, on whether this would save lines when typesetting, on whether a footnote apparatus can be shrunk by more than a certain factor in this manner, and whether the ratio of allowable joints between footnotes²⁰ to the number of footnotes exceeds a certain ratio.²¹ The criteria are configurable per apparatus or globally.

There are some footnotes where a vertical arrangement is mandatory,²² and the footnote must not be set into a hbox to start with. This is the case, for example, for footnotes containing display math. Placing a + sign before the opening brace of the footnote text will achieve that, and similarly a - sign can be used for switching in an otherwise vertically arranged footnote apparatus to horizontal arrangement.

`bigfoot` hooks into the output routine and does its accounting work before the main output routine gets a chance to get called. This work involves sorting the various contributions to a single insertion, joining together all in-paragraph footnotes into a single paragraph, measuring the resulting boxes, and gathering more material from the page in case that

²⁰ where both footnotes around the breakpoint are considered potentially horizontal material

²¹ A footnote apparatus in which there are just few horizontally arranged footnotes would appear inconsistent.

²² like footnotes containing

- list environments
- display math like

$$\sum_{n=1}^{\infty} \frac{(-1)^n}{n} = \log \frac{1}{2}$$

this produces an underfull box. Since the insertions `bigfoot` uses are unsplittable, this will often lead to an overfull box. In that case, the various footnote blocks get split to an optimum size before the real output routine gets called, and if this results in an underfull box again, more material gets called in again.

5.2 Dissecting `\@makefn`

Document classes implement the desired footnote layout with the macro `\@makefn`. This macro receives one argument, the body of the footnote. We'll now discuss several problems we want to tackle in the context of using `\@makefn` for implementing the layout prescribed by the class file.

Robust footnotes We want footnotes with robust arguments, like those of plain `TeX`, to forestall complaints when `\verb` and its catcode mongering cousins fail to work in footnotes. The trick is to have the macro argument of the `\footnote` macro not really be a macro argument, but the content of an `\hbox` or `\vbox` command, and have subsequent code do its work with `\aftergroup`, once the command finishes.

This means that we have to cut `\@makefn` into parts before and after its argument. It turns out that cutting the part before it starts processing its argument is rather easy:

```
\@makefn \iffalse \fi
```

will do that. It executes and expands `\@makefn` until it comes to the point where it would process its argument, which happens to be `\iffalse`, and then kills the rest of `\@makefn`. At least as long as the argument #1 does not happen to be in itself inside of a conditional, in which case bad things will happen. Very bad things. But a pretty thorough sampling of `\@makefn` variants on `TeX` Live did not turn up such code.

Much more problematic is getting hold of the second part of `\@makefn`. It turns out that about 95% of the variations out there in different class files will work with

```
\expandafter \iffalse \@makefn \fi
```

which looks rather similar to the above. Unfortunately, it is not quite equivalent, since in the upper code, `\@makefn` is cut into two once it has been expanded up to its macro parameter, whereas in the lower version it is cut into two before any parts of it get expanded. If any of the closing braces that follow #1 in the definition of `\@makefn` happen to belong to the argument of a macro starting before #1, they will cause spurious closing groups.

Getting the closing part at the end of the footnote without any remaining macro braces is more tricky, inefficient and error prone. One possibility is starting another instance of `\@makefntext` inside of a box to be discarded later. Then as its macro argument you use code that will repeatedly be closing opened groups until the outer group level is reached again and the box can be discarded. ϵ -`TeX`'s grouping status macros (`\currentgrouplevel` and `\currentgrouptype`) make it possible to know how to close the current group and whether it is the last one involved. After everything that has been opened has been discarded again, the remaining tokens in the input stream should form a perfect complement to the tokens that the initial `\iffalse` trick has discarded at the start of the footnote.

One other mechanism probably worth playing with is the use of alignment templates, since they provide a natural way of having `TeX` switch input contexts across groups. The best approach in that regard would seem to parse the content of the footnote within a `\noalign` group of a `\valign`, but that still suffers from the problem that no automatic discretionary are generated for explicit hyphens.

But since most of the the `\@makefntext` variants out in the field are covered with the simple variant (basically, this is the case for all definitions that do not use `#1` within a macro argument itself), `bigfoot` for now has not added any of the more complicated versions. The group discarding trick might perhaps be made available with a separate package option at a later time, if there is sufficient demand for it.

But it may be easier in most cases simply to rewrite the culprits: after all, `\@makefntext` is rarely complicated. Most notably, the `\@makefntext` of the `ltugboat` class is so ridiculously contorted that the automated analysis of it fails. (It has been replaced with an equivalent for this article.)

Using `\@makefntext` for ‘para’ footnotes is a tricky feat: the ‘para’ footnote style sets all footnotes within one continuous running paragraph, a manner of operation quite different from the original intent of `\@makefntext`. Single footnotes are first collected in horizontal mode, and at `\output` time the relevant footnotes making it to the current page are pasted together. This has several problems: for one, `\@makefntext` will set paragraph breaking parameters. We need these at the time that we assemble the footnotes into one paragraph. But `\@makefntext` also generates the footnote mark, so we need to call it for each footnote.

So even when we set `\@thefnmark`²³ equal to an empty string at footnote assembly time, the assembled footnote mark will likely take up some additional space. This is not the end of our worries: the formatting will be right for standard footnotes, but not cater for ‘para’ footnotes. If we want to have a reasonable looking turnout, here are the conditions we have to meet:

1. At the beginning of the footnote block, or if a footnote starts right after a line break, the specified formatting should be used.
2. Within the line, we shall keep the spacing between footnote mark and footnote text correct. However, most styles right-justify the footnote mark within a box of fixed size. If we keep this sort of formatting, we will end up with a large space before short footnote marks, and a small one before longer marks. Since the amount of whitespace inside of a line should not be so large as to cause unsightly white holes, nor so small to make the footnote mark confused to be a part of the preceding footnote, we want a fixed spacing before the footnote mark.

The solution to these problems is to do a few measurements: we measure the width that an empty footnote mark would cause in the footnote box (and start our assembled footnotes with a negative space compensating that), and we typeset the footnote mark once on its own with `\@makefntext`, fishing with `\unskip` and `\lastbox` for the footnote mark box and resetting it to its natural size (which will kill the particular justification prevalent in the majority of class files doing justification). The difference in box size gets recorded separately until the time that the footnote gets set, and then the interfootnote glue is calculated accordingly.²⁴

Maintaining the color stack is a particularly unwholesome field of study.²⁵

What is the color stack, anyway? `LATEX`'s `color` package provides color selection commands that will change the current text color until the end of the group, where it will be restored.

The involved macros are

`\color@begin@group` is called at the start of each ‘movable’ box: material that does not necessarily appear right away. Without color support

²³ the mark as displayed in the footnote

²⁴ A few classes work with `\parshape` or `\hangindent`, either directly or with a `list` environment, and this is also taken into consideration as far as possible.

²⁵ The main philosophy for work on the color stack has been summarized well by David Carlisle: “It’s not my fault.”

loaded, this does nothing. With color support loaded, it is usually equal to `\begin@group`.

`\color@end@group` is the corresponding macro at the end of ‘movable’ boxes. Any color restoration initiated with `\aftergroup` in the box will happen right here, still within the scope of the box, instead of outside where it would not move with the box.

`\set@color` will be called for setting the current color. It will also use `\aftergroup` in order to insert a call to `\reset@color` when the group ends.

`\reset@color` will restore the current color to what it was before the current group.

How will the color be restored? We have two different models:

dvips restores colors by making use of a color stack: `dvips` can ‘push’ a new color onto the stack, and pop the previous color back. Consequently, `\reset@color` inserts a special that tells `dvips` to pop the stack once.

pdftex instead restores colors by reinstating the color stored in `\current@color` after closing the group.²⁶

It is clear that the `pdftex` model is insufficient to even keep the color of the main text across page breaks, since on the next page there is no special after the page break that could switch back to the text color after the page footer²⁷ from the last page and headers from the current page have been placed with a default color.²⁸

But in the context of footnotes, the problem is severely exacerbated: a footnote can be broken right in the middle of a sequence of color changes. The technically sound solution would be to switch to a different color stack for each footnote block. Since `dvips` does not offer multiple color stacks (and `pdftex` does not even offer a single one), we have to revert to trickery.

At each color change, the complete state of the color stack gets recorded in a mark. When the footnote is broken, we use the information in the mark in order to unwind the color stack to the state on the page before the footnote was entered. When the footnote is continued on the next page, the unwound color stack is reinstated again. Whenever `\color@begin@group` is called, the whole recording and restoration business is stopped (since a new

²⁶ Of course this means that if we are at the end of a movable box, the restored color will be that at the time the box was assembled, not at the time it was used.

²⁷ and footnotes

²⁸ Heiko Oberdiek’s `pdfcolmk` package tries to deal with that particular problem.

context has been started), the record of the color stack essentially restored to empty, and only resumed when the corresponding group has ended.

In order to keep these proceedings fit for consumption by the general public, the reader is referred to the actual code for further details.

6 Outlook

At the time this article was written, quite a few tasks remained to be done. Further improvements in the footnote breaking decisions and their scoring metrics are needed. Flushing footnotes out in the middle of the page for short successive works would be nice. Amending footnotes with marginals (including line numbers) in a manner consistent with the main text would seem desirable. Additional footnote arrangements apart from the existing basic two styles should be easily implementable on top of the general scoring and breaking mechanisms.

7 Conclusion

It is hoped and expected that this bundle will become a basic building block for critical typesetting applications. While there are other packages available for that purpose, `bigfoot` (with its companions) offers the following important features:

- It is completely layout-neutral: while most solutions for critical typesetting are provided in the form of document classes, `bigfoot` does not make layout decisions but instead just uses the layout provided by a base class.
- Footnote arrangement and balancing is vastly superior to and more flexible than any of the other available solutions.
- Color works.
- The interfaces for creating new functionality focused around footnotes are reasonably simple.

At the time this article was written, not all interfaces have been cast into stone. However, `bigfoot` can be mostly used as an upwards-compatible drop-in replacement of `manyfoot`.

All that remains is to profusely apologize for the quite inappropriate use²⁹ of footnotes in this article for illustrative purposes.

References

- [1] <http://sarovar.org/projects/bigfoot>
(developer site and CVS instructions)
- [2] `CTAN:macros/latex/contrib/bigfoot`
(released packages)

²⁹ or rather abuse

μονο $\vec{2}$ πολυ: Java-based Conversion of Monotonic to Polytonic Greek

Johannis Likos

ICT Consulting
Rusthollarintie 13 E 35
Helsinki 00910
Finland
likosjo@yahoo.com

Abstract

This paper presents a successfully tested method for the automatic conversion of monotonic modern Greek texts into polytonic texts, applicable on any platform. The method consists of combining various freely available technologies, which have much better results than current commercially available solutions. The aim of this presentation is to introduce a way of applying this method, in order to convert thousands of digitally available single-accented modern Greek pages into attractive artworks with multi-accented contents, which can be easily transferred either to the Web or a T_EX-friendly printer. We will discuss the preparatory and postprocessing efforts, as well as the editing of syntax rulesets, which determine the quality of the results. These rulesets are embedded in extendable tables, functioning as flat databases.

1 Introduction

During the past centuries, Greek and Hellenic scholars have introduced and refined polytonism (multiple accenting) in the written word for the precise pronunciation of ancient Greek. Since spoken modern Greek is comparatively less complicated, the Greek government has officially replaced polytonism by monotonism (single accenting) for purposes of simplification, especially in the educational system. Also, Greek authors commonly use monotonism, since it is so much simpler to produce.

Classical, polytonic, Greek has three accents (acute, grave, and circumflex) and two *breathings* (rough and smooth—equivalent to an initial ‘h’ and lack thereof). Accents are lexically marked, but can change based on other factors, such as *clitics* (small, unstressed words that lean on another word to form a *prosodic word*—a single word for accent placement). In addition, two other symbols were used: diaeresis (to indicate two vowels that are not a diphthong) and iota subscript (a small iota that was once part of a diphthong but subsequently became silent).

Monotonic Greek retains only the acute accent, which was usually, though not always, the same as the classical acute. To make a graphic break with the past, the new acute accent was written as a new *tonos* glyph, a dot or a nearly vertical wedge, although this was officially replaced by a regular acute in 1986.

So, why bother with the complexities of polytonism? The benefits are increased manuscript readability and, even more important, reducing ambiguity. Despite the simplification efforts and mandates, the trend nowadays is *back to the roots*, namely to polytonism. More and more publishers appreciate, in addition to the content, the public impression of the quality of the printed work.

This paper discusses an innovative and flexible solution to polytonism with an open architecture, enabling the automatic multiple accenting of existing monotonic Greek digital documents.

2 Terminology

In this article, we will use the terms *polytonism* and *multiple accenting* interchangeably to mean the extensive usage of *spiritus lenis*, *spiritus asper*, *iota subscript*, *acute*, *gravis* and *circumflex*. Similarly, we use the terms *monotonism* and *single accenting* to mean the usage of simplified accenting rules in Modern Greek documents.

3 Historic Linguistic Development

During the last four decades the printed Greek word has undergone both minor and radical changes. Elementary school text during the late 1960s and early 1970s made Purified Greek (καθαρεύουσα) imperative, by strict government law of the time. The mid-1970s saw a chaotic transition period from Purified Greek to Modern Greek (δημοτική) with simplified

grammar, where some publications were printed with multiple accenting, some with single accenting and even some without any accenting at all!

Even after the government officially settled on monotonism in the early 1980s, Greek publishers were not able to switch immediately to the monotonic system. During the last decade, many computerized solutions have been invented for assistance in typing monotonic Greek. Today, there is a trend toward a mixture of simplified grammar with multiple accenting, decorated with Ancient Greek phrases. See table 1.

4 Polytonic Tools

There are two programs for Microsoft Word users, namely *ΤΟΝΙΣΜΟΣ* by DATA-SOFT and *ΑΥΤΟΜΑΤΟΣ ΠΟΛΥΤΟΝΙΣΤΗΣ* (academic and professional version) by MATZENTA. A third is the experimental *μονοῶπολυ*, an open source project, which is the subject of this discussion.

These solutions are independent. The major difference between the commercial and open source programs is the control of the intelligence, such as logic, rule sets and integrated databases. In the case of the commercial solutions, users depend on the software houses; in the open source case, users depend on their own abilities. See table 2.

There is no absolutely perfect tool for polytonism, so the ultimate choice is of course up to users themselves.

5 Open Source Concept

μονοῶπολυ implements a modular mechanism for multiple accenting of single-accented Greek documents. See figure 1.

5.1 Architecture

The *μονοῶπολυ* architecture consists of (figure 2):

- methods: DocReader, DocWriter, DBParser, Converter
- configuration file: *.cfg
- flat database: *.xml
- document type definition: *.dtd
- optional spreadsheet: *.csv, *.xls

5.2 Configuration

The plain text configuration file defines (gray arrows in fig. 2) the necessary filenames and pathnames. The *μονοῶπολυ* components read this during initialization to determine where to find the input files and where to write the output files (by default, the current working directory).

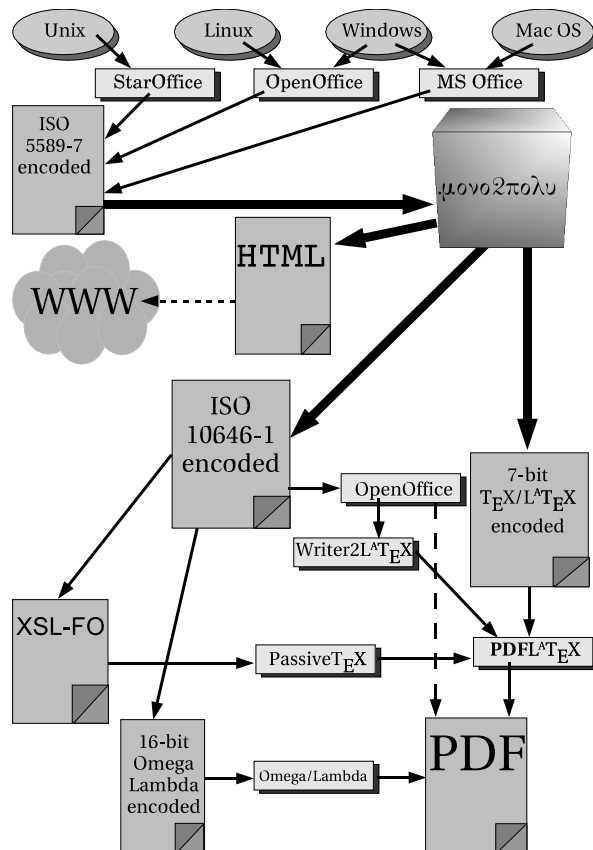


Figure 1: Overview of the overall multiple accenting concept, which involves many external tools.

5.3 Database Connectivity

The dotted arrows in the architecture figure (fig. 2) show the connection between a CSV spreadsheet, a Document Type Definition (DTD), the actual XML flat database, and the database parser.

5.4 Input Prerequisites

During the conversion process, invisible special control codes for formatting features (superscript, bold, etc.) make it difficult to coherently step through paragraphs, sentences and words. Therefore, plain text files serve best for polytonism input.

The DocReader component of *μονοῶπολυ* expects the source document to be in the ISO 5589-7 encoding, and to be written according to the Modern Greek grammar, especially regarding the monotonic accenting rules.

Assistance for monotonic accenting while typing Modern Greek documents is provided by Microsoft's commercially bundled spellchecker, or any downloadable Open Source spellchecker.

Table 1: Selected examples of Greek publications from the last four decades.

Date	Publisher	Author	Subject/Title	Language	Remarks	Example
1968	ΕΚΔΟΣΕΙΣ Α. ΚΑΡΑΒΙΑ	S. TIMOSHENKO D.H. YOUNG	Ἀντοχή τῶν Ὑλικῶν	purified	polytonic translation including gravis	
1971	Ο.Ε.Δ.Β.	Μ. ΚΑΤΣΙΚΑΣ	ΓΕΩΓΡΑΦΙΑ ΣΤ' ΔΗΜΟΤΙΚΟΥ	purified	no dative at all and iota subscript seldom used	νὰ ἐπιδιωχθῆ
1978	ΙΔΡΥΜΑ ΕΥΓΕΝΙΔΟΥ	Ι. ΧΑΣΤΑΣ	βιβλ. Τεχν. κ. Ἐπαγγ. Λυκείου	modern	polytonic without gravis	
1980	self-published	Ρ. ΓΡΑΙΚΟΥΣΗΣ	ΣΤΟΙΧΕΙΑ ΜΗΧΑΝΩΝ	modern	acute used instead of gravis in polytonic text, some feminine singular genitive in purified version	τῆς κλίσεως
1982	self-published	Ε.Σ.Μ.Α.	ΠΡΑΚΤΙΚΑ ΕΘΝ. ΑΕΡΟΠΟΡΙΚΟΥ ΣΥΝΕΔΡΙΟΥ	modern	monotonic typewriter text style	
1989	INTERBOOKS	ΓΡ. ΣΦΑΚΙΑΝΟΣ	ΕΜΠΟΡΙΚΗ ΑΛΛΗΛΟΓΡΑΦΙΑ	modern	monotonic with neutral accent and acute	εκατός από τόν
1998	ΠΑΡΑΤΗΡΗΤΗΣ	Α. ΣΥΤΡΟΠΟΥΛΟΣ	ΛΑΤΕΧ	modern	acute only accent type in monotonic text	τό τόπι
2003	Α. ΦΩΤΙΕΡΗΣ	several authors	ἡ λέξιη	modern	polytonic without gravis	καί ἀπ' τό

Table 2: Comparison of polytonic tools.

	ΤΟΝΙΣΜΟΣ	ΑΥΤΟΜΑΤΟΣ ΠΟΛΥΤΟΝΙΣΤΗΣ	μονοῤπολυ
Greek language support			
Ancient	yes	yes	later
Hellenistic	yes	yes	later
Byzantine	no	no	later
Church/Biblical/NT	yes	yes	later
Purified	yes	yes	later
Modern	yes	yes	yes
Mixed (Ancient and Modern)	selectable	selectable	fixed
Editing assistance	no	yes	no
Database	fixed (binary)	fixed (binary), editable exception list	editable (XML)
Manual corrections		interactively	post-processed
Automatic hyphenation	unknown	yes	external task
Protection	hardlock	ID	no
File formats			
Input	Word (Win, Mac)	Word (Win, Mac) other Greek formats	ISO 5589-7 encoded ASCII on any platform
Output	Word (Win, Mac)	Word (Win, Mac) other Greek formats, HTML	ASCII (ISO 10646-1 encoded), HTML (ISO 10646-1 encoded)
Unicode support	yes	yes	yes
TeX-specific filters	TeXto737.lex, 737toTeX.lex	no	WriterῤLaTeX
Requirements			
Platforms			
Microsoft Windows	95, 98, ME, NT, 2000, XP	Microsoft Word	JDK 1.4
Apple Macintosh	no	no	95, 98, ME
Linux	no	no	NT, 2000, XP
Unix	no	no	Mac OS 9, Mac OS X
Availability	immediately	immediately	Mandrake, Red Hat, SuSE
Distribution	purchased license	purchased license	AIX, HP-UX, Sinix, Solaris
			under development
			open source

5.5 Converter

The bold arrows in the architecture figure (fig. 2) show the data exchange between the internal components, the document reader, the database parser and the document writer to the converter. The conversion process does not include grammar analysis, since μονοῤπολυ expects that monotonic proof reading has been done previously, with other tools.

6 External Interfaces

The output from μονοῤπολυ (DocWriter method) is in the ISO 10646-1 encoding, in various formats,

which are then post-processed. The dashed arrows in fig. 2 show the relationship between the external files.

6.1 Web Usage

For background, these web pages discuss polytonic Greek text and Unicode¹ (UTF) fonts:

¹ Concerning missing Greek characters and other linguistic limitations in Unicode, see *Guidelines and Suggested Amendments to the Greek Unicode Tables* by Yannis Haralambous at the 21st International Unicode Conference

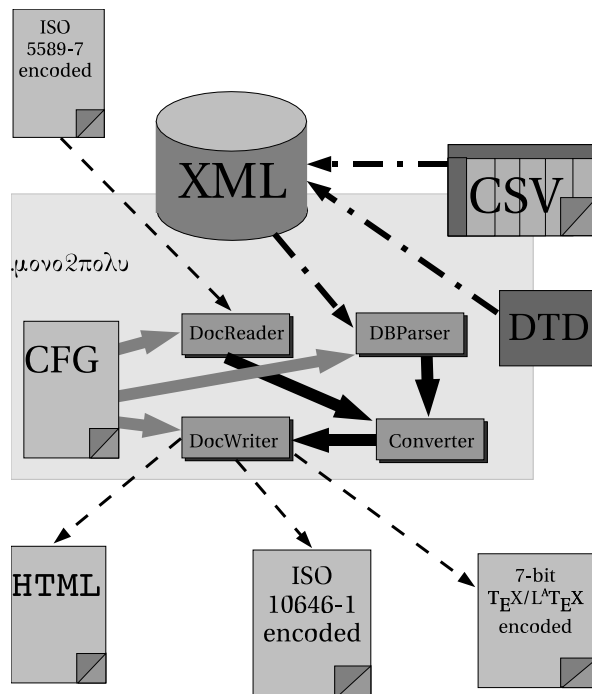


Figure 2: Overview of internal architecture, with external interfaces to existing standards.

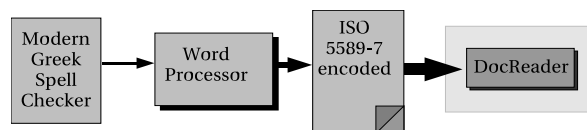


Figure 3: External creation of monotonic Greek with any word processor (e.g., OpenOffice) using a separate spellchecker (e.g., elspell).

- <http://www.ellopos.net/elpenor/lessons/lesson2.asp>
- <http://www.stoa.org/unicode/>
- <http://www.mythfolklore.net/aesopica/aphthonius/1.htm>

on May 2002 in Dublin, Ireland (<http://omega.enstb.org/yannis/pdf/amendments2.pdf>).

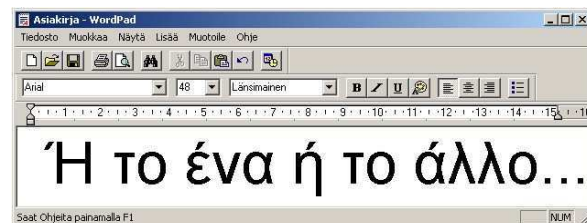


Figure 4: Example of original monotonic input.

Using the direct HTML polytonic output from *μονοῤπολυ* requires that the layout of the web page be done in advance, since manually editing the numeric Unicode codes in the *.html file is impractical (see figure 5). Dynamic web pages created through CGI scripts, PHP, etc. have not yet been tested.

```
<font
FACE="Arial Unicode MS"
SIZE="36">
&#7978;&#32;
&#964;&#8056;&#32;
&#7957;&#957;&#945;&#32;
&#7970;&#32;
&#964;&#8056;&#32;
&#7940;&#955;&#955;&#959;
&#46;&#46;&#46;
</font>
```

Figure 5: Example of polytonic HTML output.

6.2 OpenOffice Usage

The ISO 10646-1 encoded polytonic output (fig. 6) from *μονοῤπολυ* could be inserted into the OpenOffice Writer software, since the newest version can directly output polytonic Greek .pdf files. Unfortunately, the quality of the result leaves much to be desired. Better results can be produced by converting from Writer to LḂTḂX and doing further processing in the LḂTḂX environment.



Figure 6: Example of polytonic output.

6.3 (LḂ)TḂX Usage

The most likely scenario for (LḂ)TḂX users is using the Greek babel package, and adding the *μονοῤπολυ* 7-bit polytonic output text into the source .tex file. See figures 7 and 8.

The 7-bit output from *μονοῤπολυ* could presumably also be inserted into .fo files, and processed through PassiveTḂX, but this has not yet been tested. Likewise, the ISO 10646-1 output could presumably be processed directly with Ω/LḂ, but this has not been tested, either.

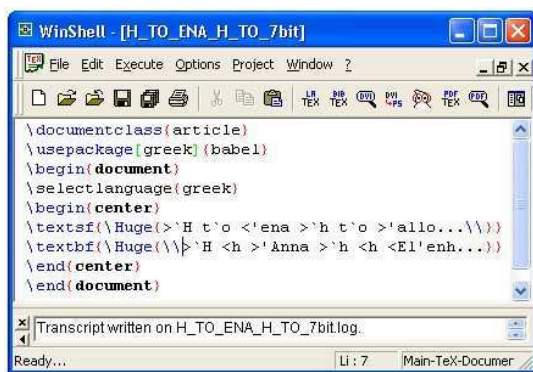


Figure 7: Example of polytonic T_EX output, either from *μονοῶπολυ* or Writer2L^AT_EX.

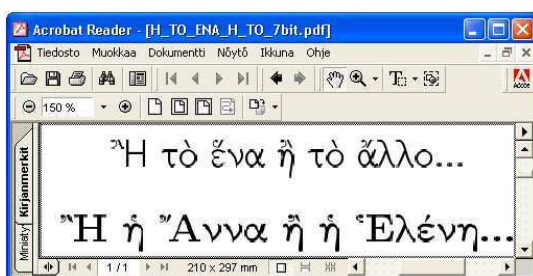


Figure 8: Polytonic PDF output from T_EX.

7 Technologies Used in *μονοῶπολυ*

After some evaluation, we chose to focus on Java, Unicode and XML, due to their flexibility in processing non-Latin strings, obviously a critical requirement of *μονοῶπολυ*.

7.1 Programming Language

Two major reasons for choosing Java (J2SE) as the implementation language of *μονοῶπολυ* were the capabilities for handling XML and Unicode through widely-available and well-documented libraries. The Java SDK provides extremely useful internationalization features, with the ability to easily manipulate string values and files containing wide characters.

In order to concentrate on *μονοῶπολυ*'s essential features, no graphical user interface has been designed.

7.2 Character Set

The choice of Unicode/ISO 10646-1 for the character set should be clear. It combines monotonic and polytonic Greek letters, is known worldwide and standardized on most platforms, and contains most (though not all) variations of Greek vowels and

consonants, in the Greek and the Greek Extended tables.²

For further information on writing polytonic Greek text using Unicode, see <http://www.stoa.org/unicode/>.

7.3 Text Parsing Libraries

Most helpful for the parsing of XML-based database entries are the SAX and DOM Java libraries.

The following Java source code, taken from the *μονοῶπολυ* class *DBparse*, serves to demonstrate usage of SAX and DOM. The code counts and then outputs the total amount of all available entries in the XML database file.

```
import java.io.*;
import org.xml.sax.SAXException;
import org.xml.sax.SAXParseException;
import javax.xml.parsers.DocumentBuilder;
import javax.xml.parsers.DocumentBuilderFactory;
import javax.xml.parsers.FactoryConfigurationError;
import javax.xml.parsers.ParserConfigurationException;
import org.w3c.dom.*;

public class DBparse{
    static Document document;
    String warn="No XML database filename given...";
    public static void main(String param[]){
        if (param.length!=1){
            System.out.println(warn);
            System.exit(1);}
        File mydbfile=new File(param[0]);
        boolean load=mydbfile.canRead();
        if (load){
            try{
                DocumentBuilderFactory fct
                    = DocumentBuilderFactory.newInstance();
                DocumentBuilder builder
                    = fct.newDocumentBuilder();
                document = builder.parse(mydbfile);}
            catch (SAXParseException error){
                System.out.println("\nParse error at line: "
                    + error.getLineNumber() + " in file: "
                    + error.getSystemId());
                System.out.println("\n" + error.getMessage() );}
            catch (ParserConfigurationException pce)
                {pce.printStackTrace();}
            catch (IOException ioe){ioe.printStackTrace();}
            catch (Throwable t){t.printStackTrace();}}
        else{System.out.println("XML database missing!");}
        String mytag='\u03C3'+"";
        NodeList taglist=document.getElementsByTagName(mytag);
        int amount=taglist.getLength();
        System.out.println("amount of entries:\n" + amount );}}
```

Notice particularly the fourth-last line, where *mytag* is assigned '\u03C3', namely the character σ , used as the search string.

8 Database Structure

The XML standard from the W3C has proven to be a simpler choice for storing either monotonic or polytonic Unicode text than the alternatives, such as spreadsheets or even SQL databases. The quality of the final polytonic result depends on the precision of the XML content, where ambiguities have to

² <http://www.unicode.org/versions/Unicode4.0.0/ch07.pdf>

be marked with special symbols for manual post-processing.

Currently, the entries of the basic database consist of tags with parameters and values. The tag name indicates the type of the expression: a single character, a prefix, a suffix, a substring, a word or a chain of words. The five parameters are as follows:

1. The monotonic ISO 5589-7 encoded source expression to be converted.
2. The Unicode output text.
3. A 7-bit output text for (L^A)T_EX usage with the Greek `babel` package.
4. The equivalent numeric value according to the *Extended Greek* Unicode table for HTML usage.
5. An explanatory comment or example, in case of ambiguities or linguistic conflicts.

Here, I have built on the work of prior Greek T_EX packages, such as GreekT_EX (K. Dryllarakis), Scholar T_EX (Y. Haralambous), and `greektex` (Y. Moschovakis and G. Spiliotis), for techniques of using the iota subscript, breathings and accents in 7-bit transliterated `.tex` source files.

In the following examples, note carefully the different bases used: '074 is octal, #8172 is decimal and '03D1' is hexadecimal.

8.1 Data Type Definition

The required basic *Data Type Definition* is currently located in the experimental namespace `xmlns:β = http://koti.welho.com/ilikos/TeX/LaTeX/mono2poly/mono2poly.dtd`. It contains the following information:

```
<!ELEMENT β (σ+)>
<!ELEMENT σ (#PCDATA)>
<!ATTLIST σ
  μ CDATA #REQUIRED
  π CDATA #REQUIRED
  τ CDATA #REQUIRED
  δ CDATA #REQUIRED
  ξ CDATA #REQUIRED>
```

Thus, we have one element, called β (βάση δεδομένων = *database*). It contains multiple element sets, called σ (στοιχεία συλλαβῆς = *syllable data*). Each element set has, at present, five attributes, namely μ for monotonic expressions, π for polytonic expressions, τ for 7-bit (L^A)T_EX code, δ for HTML code, and finally ξ for comments.

The DTD can be overridden by a local `.dtd` file, which must be specified in the header of the `.xml` database file; for example:

```
<!DOCTYPE β SYSTEM "my_own_mono2poly.dtd">
```

Both the `.dtd` and `.xml` must reside in the same directory.

8.2 Data Entries

Here is an example database entry, showing the only Greek capital consonant with spiritus asper:

```
<σ
  μ="P"
  π="Ψ"
  τ="\char'074 R"
  δ="#8172;"
  ξ="Ψόδοϛ"
></σ>
```

The *slash* symbol indicates the closing element tags in XML, while the *backslash* symbol is used for (L^A)T_EX commands. Both appear in the `.xml` database file.

8.3 Header and Body

Although not explicitly documented, exotic characters may be used in `.dtd` and `.xml` files as long as the appropriate encoding is declared:

```
<?xml version="1.0" encoding="UTF-16"?>
```

The header should include other information as well. Schematically:

```
<?xml version="1.0" encoding="UTF-16"?>
<!DOCTYPE β SYSTEM "mono2poly.dtd">
<!-- author: ... -->
<!-- affiliation: ... -->
<!-- creation date: ... -->
...
<!-- notes: ... -->
<β>
  <σ μ="..." π="..." τ="..." δ="..." ξ="..."></σ>
  ...
  <σ μ="..." π="..." τ="..." δ="..." ξ="..."></σ>
</β>
```

For quality assurance, after database creation and after each update a validation and verification test should be run, to detect XML syntax errors and linguistic content mistakes.

This concept of the database as a lookup/mapping table allows differentiating between initial and intermediate consonants. For example:

```
β (03D0) ↔ β (03B2)
θ (03D1) ↔ θ (03B8)
ρ (03F1) ↔ ρ (03C1)
φ (03D5) ↔ φ (03C6)
```

Therefore, by updating the XML file, post-processing may be reduced. Experienced linguists may wish to use different tools for the correcting and the updating of the flat database. Rows with multiple columns from spreadsheets can be inserted directly into XML data files, as long as the columns are sorted in the expected order.

8.4 Expression Types

In each database entry, there is one source expression, at least three target expressions, and possibly one explanation. The ISO 5589-7 encoded source expression and the first ISO 10646-1 encoded target expression may be a:

- single uppercase or lowercase character with or without spiritus and/or accent
- partial word, such as prefix, intermediate syllable, suffix
- complete word
- chain of combined words
- combination of partial word pairs, such as a suffix followed by a prefix
- mixture of complete and partial words, such as a complete word followed by a prefix or a suffix, followed by a complete word

The rest of the target expressions represent the same information as the first in other output formats, namely for 7-bit Greek (L^A)_{TEX} and HTML as well. The intelligence of the μνο2πολυ system currently lies in the database, so while creating and editing entries, it is crucial to write them correctly.

8.5 Editing Tools

One of the most powerful Unicode editors is the Java-based Simredo 3.x by Cleve Lendon, which has a configurable keyboard layout, and is thus suitable for this sort of task. The latest version of Simredo, 3.4 at this writing, can be downloaded from <http://www4.vc-net.ne.jp/~klivo/sim/simeng.htm>, and installed on any platform supporting JDK 1.4.1 from Sun. Simredo can be started by typing `java Simredo3` or perhaps `java -jar Simredo3.jar` in the shell window (Linux) or in the command window (Windows). Unicode/XML with Simredo has been successfully tested on Windows XP and on SuSE Linux 8.1 Professional Edition.

The author would be happy to assist in the preparation of a polytonic Greek keymap file (.kmp) for Simredo, but the manual may prove sufficient. The creation of such a keymap file is easily done by simply writing one line for each key sequence definition. For instance, given the sequence `2keys;A^A` using the desired Unicode character, or the equivalent sequence `2keys;A\u1F0D` with the big endian hexadecimal value, one can produce an uppercase *alpha with spiritus asper and acute accent* by pressing the `[;]` and `[A]` keys simultaneously. According to the Simredo manual, other auxiliary keys such as `Alt` can be combined with vowel keys, but not `Ctrl`.

Some other Unicode editors:

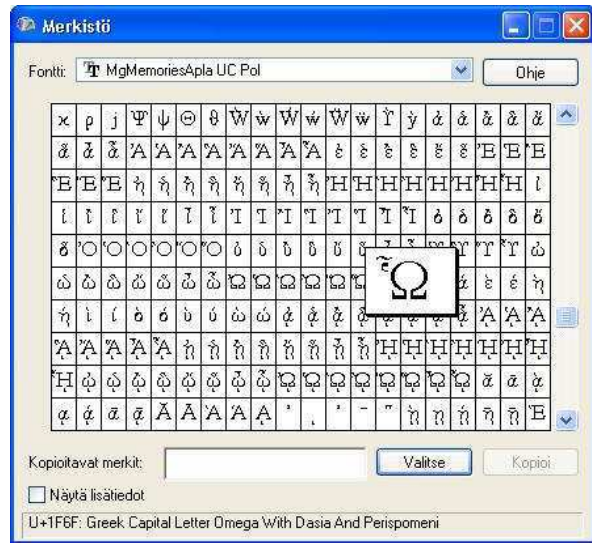


Figure 9: Another useful tool is a character mapping table like this accessory on Windows XP, which displays the shape and the 16-bit big endian hexadecimal code of the selected character.

- For Windows: http://www.alanwood.net/unicode/utilities_editors.html.
- For Linux: <http://www.unicodecharacter.com/unicode/editors.html>.
- For Mac OS: <http://free.abracode.com/sue/>.

Unfortunately, XMLwriter and friends neither support configurable keyboard layouts nor display 16-bit Unicode.

8.6 Polytonic Keyboard Drivers

Instant interactive multi-accenting while editing Greek documents is available either through plug-ins for some Windows applications, such as SC Unipad (<http://www.unipad.org/main/>) and Antioch (<http://www.users.dircon.co.uk/~hancock/antioch.htm>), or with the help of editable keyboard mapping tables, such as the Simredo Java program described above. Regrettably, the *Hellenic Linux User Group* (HELL.U.G., <http://www.hellug.gr> and <http://www.linux.gr>) has no recommendations for polytonic Greek keyboard support.

Whatever polytonic keyboard driver has been installed and activated may be useful for new documents, but does not much help the author who is not familiar with the complicated rules of polytonism!

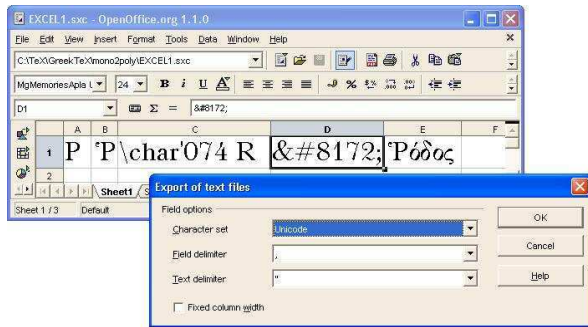


Figure 10: Using a spreadsheet to produce a long extendable list with five columns, which then can be saved as a .csv file. Be careful with the parametrization!

8.7 Auxiliary Tables

Preparation and periodic updates of auxiliary tables can of course be done with any software supporting Unicode. Spreadsheets have the advantage of putting the entries into cells row-by-row and thus organizing the parameters by column. This may prove easier than directly writing the XML file. See figure 10.

A row in such a .csv file looks like this:

"P", "P", "\char'074 R", "Ῥ ;", "Pόδος"

Of course it then must be re-shaped with element and attribute tags to make an XML-syntax database entry.

8.8 Viewing Tools

Users without any programming knowledge may find it useful to open and inspect the header and the body of the XML database before using it in polytonic documents. Here is a procedure for doing that.

First, set the Unicode font in the preferences of the desired browser (fig. 11). These days, most browsers support this, including Internet Explorer, Konqueror, Netscape Navigator and Opera.

Then, select Unicode UTF-16 as the default encoding (fig. 12). The browser can now detect syntax errors, giving immediate feedback (fig. 13).

8.9 Priorization of Database Entries

Polytonic exceptions (e.g., οὔτε and ὄστε without circumflex) and especially ambiguities (e.g., που → ποὺ or πού, πού → ποῦ; πως → πώς or πῶς, πῶς → πῶς) have the highest priority in the database, then the special expressions, while the simple, casual and obvious accented syllables or particles have the lowest priority. In order to avoid mis-accented and

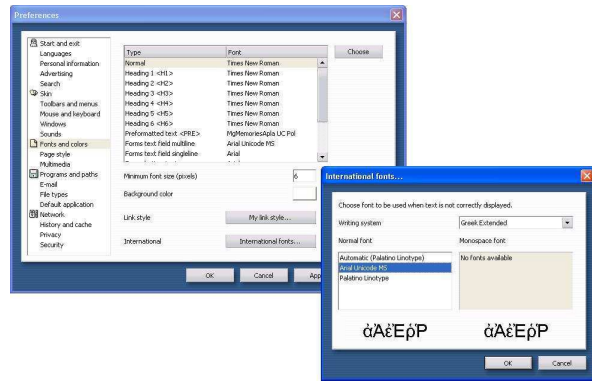


Figure 11: Choosing the Unicode font for viewing in a browser.

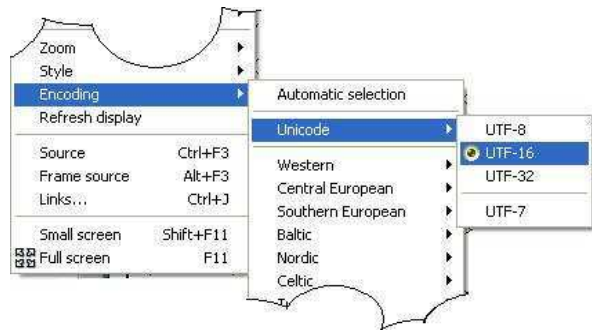


Figure 12: Selecting UTF-16 for the default encoding.

mis-spirited syllables as much as possible, entries must be in the appropriate order.

For example, table 3 shows lexical rules defining eight variations of the Greek interrogative pronoun τί (= “which”) as a single monotonic expression:

- with and without neutral accent
- with and without Greek question mark
- standalone
- leading word in the sentence
- intermediate word in the sentence
- trailing word in the sentence

Database entries like these are needed to account for the variations shown in table 4. As a rough analogy in English, it is as if table 3 shows variations on “I”: Initial position (“I went to the store”); after



Figure 13: Example error message from browser.

Table 3: Eight variations of τί as a monotonic expression.

```
<!-- ἐρωτηματικές ἀντωνυμίες -->
<σ μ="Τι;" π="Τί;" τ="Τ'ι;" δ="ᾠ932;ᾠ8055;ᾠ894;" ζ="Τί;" > </σ>
<σ μ="Τι " π="Τί " τ="Τ'ι " δ="ᾠ932;ᾠ8055;ᾠ32;" ζ="Τί λές;" > </σ>
<σ μ=" τι;" π=" τί;" τ=" τ'ι;" δ="ᾠ32;ᾠ964;ᾠ8055;ᾠ894;" ζ="Κι' ἔγραψε τί;" > </σ>
<σ μ=" τι " π=" τί " τ=" τ'ι " δ="ᾠ32;ᾠ964;ᾠ8055;ᾠ32;" ζ="Καί τί ἔγραψε;" > </σ>
<σ μ="Τί;" π="Τί;" τ="Τ'ι;" δ="ᾠ932;ᾠ8055;ᾠ894;" ζ="Τί;" > </σ>
<σ μ="Τί " π="Τί " τ="Τ'ι " δ="ᾠ932;ᾠ8055;ᾠ32;" ζ="Τί λές;" > </σ>
<σ μ=" τί;" π=" τί;" τ=" τ'ι;" δ="ᾠ32;ᾠ964;ᾠ8055;ᾠ894;" ζ="Κι' ἔγραψε τί;" > </σ>
<σ μ=" τί " π=" τί " τ=" τ'ι " δ="ᾠ32;ᾠ964;ᾠ8055;ᾠ32;" ζ="Καί τί ἔγραψε;" > </σ>
```

Table 4: Similar but unrelated syllables treated differently.

Position	Syllable	Polytonic examples
leading	τι-, Τι-	τιμή, Τισσαφέρνης
intermediate	-τι-	ἀδυνάτισμα
trailing	-τι	κάτι, πράγματι
leading	τι-, Τι-	τίποτα, Τίγρης
intermediate	-τι-	ἐκτίμηση
trailing	-τί	γατί, ψωμί, τυρί

a verb (“What do I know?”); etc., and then table 4 shows that “I” isn’t always capitalized: “It looks good” vs. “Can you see it?”

The above does not cover all cases related to τί. The monotonic text γιατί may be accented in two different ways in polytonic Greek (namely, γιατί and γιατί); additional entries would be required to handle this.

9 Polytonic Printing Press

The author has found several Greek newspapers, magazines, and books, including university presses, using polytonic Greek:

- daily newspapers:
 - Η ΚΑΘΗΜΕΡΙΝΗ
 - ΕΣΤΙΑ
- monthly magazines:
 - ΝΕΜΕCΙC
- book publishers:
 - ΕΚΔΟΤΙΚΗ ΑΘΗΝΩΝ
 - ΚΥΡΙΑΚΙΔΗΣ
 - ΓΕΩΓΡΙΑΔΗΣ
 - ΠΑΠΑΝΙΚΟΛΑΟΥ
 - ΙΝΔΙΚΤΟΣ
 - ΓΝΩCΗ
 - ΚΑΛΟΦΩΛΙΑC

- academic, polytechnic and university presses:
 - Academy of Athens
 - Polytechnics of Athens
 - University Publications of Crete
 - University of Ioannina
 - Democritian University of Thrace
- private educational organizations:
 - ΚΟΡΕΑΚΟ
 - ΜΩΡΑΪΤΗ
- others:
 - Hellenic Parliament
 - military press
 - Orthodox Church

10 Testing

The following testing procedure was used for μονοῦπολυ development. The author worked on SuSE Linux and Windows XP, but any platform (Linux, Unix, Mac or Windows) should work as well, as long the JDK is installed.

1. Visit any web site with rich Modern Greek content, for example, news sources such as <http://www.pathfinder.gr>.
2. Open a new document with a word processor supporting spell checking of monotonic Greek.
3. Copy a long excerpt of continuous text from the web site.
4. Paste the selected and copied text into the word processor window.
5. Correct any misspelled words, but do not use any style or font effects.
6. Save the document as plain ISO 5589-7 encoded text file.
7. Process the document with μονοῦπολυ, as a Java application from a console window.
8. Take the 7-bit T_EX result and add it to the L^AT_EX template file in your favourite environment (LyX, Kile, MiK_TE_X, etc.).

9. Produce a .ps or a .pdf file and check the final result with GSview or some other reader.

The results improve as the database is enriched. However, some manual editing is inevitable, depending on the complexity of the document to be multi-accented, because authors may mix Ancient Greek phrases into Modern Greek sentences.

11 Future Developments

One important improvement would be to relocate some of the intelligence to external script files, for defining and modifying the polytonic grammar rule sets.

Another avenue is to integrate $\mu\omicron\nu\sigma\tilde{\rho}\omicron\lambda\nu$ with the source and data files of the open source Writer2 \LaTeX project (by Henrik Just, <http://www.hj-gym.dk/~hj/writer2latex/>). That would provide a reverse conversion, from UTF-16BE/LE encoded Greek documents into 7-bit (\LaTeX) \TeX .

12 Related links

Here we list some further readings on the complexity of Greek multiple accenting and related subjects. First, these articles (mainly written in Greek) on the importance of the spiritus lenis and especially of the spiritus asper:

- http://www.typos.com.cy/nqcontent.cfm?a_id=4681
- <http://www.kairatos.com.gr/polytoniko.htm>
- <http://www.krassanakis.gr/tonos.htm>
- http://www.mathisis.com/nqcontent.cfm?a_id=1767

Further general sources are the following:

- Ministry of National Education and Religion Affairs — http://www.ypepth.gr/en_ec_home.htm
- Institute for Language and Speech Processing — <http://www.ilsp.gr>
- <http://www.ekivolos.gr>

References

- Blomqvist, J., Toivanen, A., *Johdatus Uuden testamentin Kreikkaan*. Yliopistopaino: Helsinki (1993), 15–25.
- Bornemann, E., Risch, E., *Griechische Grammatik*. Verlag Moritz Diesterweg: Frankfurt a./M. (1978), 1–25.
- Deitsch, A., Czarnecki, D., *Java Internationalization*. O'Reilly: Beijing, Cambridge, Köln, Paris, Sebastopol, Taipei, Tokyo (2001), 171–198.
- Harold, E.R., *Java I/O*. O'Reilly: Beijing, Cambridge, Köln, Paris, Sebastopol, Taipei, Tokyo (1999), 387–413
- Kaegi, A., Bruhn, E., *Griechische Schulgrammatik*. Verlag Weidmann: Zürich, Hildesheim (1989), 2–10.
- Dr.Maier, F., Weiss, M., Zeller, A., *OPTANON Grammatik I+II*. Bayerischer Schulbuch-Verlag: München, C.C. Buchners Verlag, Salzburg (1979), 1–12.
- Οικονόμου, Μ.Χ., *ΓΡΑΜΜΑΤΙΚΗ ΤΗΣ ΑΡΧΑΙΑΣ ΕΛΛΗΝΙΚΗΣ. ΙΝΣΤΙΤΟΥΤΟ ΝΕΟΕΛΛΗΝΙΚΩΝ ΣΠΟΥΔΩΝ: Θεσσαλονίκη* (1993), 12–21.
- Pakarinen, W., *Kreikan kielioppi*. Suomalaisen Kirjallisuuden Seura: Helsinki (1993), 1–13.
- Seeboerger-Wechselbaum, M., *Java/XML. DAS bhv TASCHENBUCH*. verlag moderne industrie: Bonn Landsberg (2002), 181–208.
- Syropoulos, A., Tsolomitis, A., Sofroniou, N., *Digital Typography Using \LaTeX* . Springer Professional Computing, Springer-Verlag: Berlin, Heidelberg, New York (2003), 315–319.
- Τσάρτζανος, Α., *ΓΡΑΜΜΑΤΙΚΗ ΤΗΣ ΑΡΧΑΙΑΣ ΕΛΛΗΝΙΚΗΣ ΓΛΩΣΣΗΣ. ΕΚΔΟΤΙΚΟΣ ΟΙΚΟΣ ΑΔΕΛΦΩΝ ΚΥΡΙΑΚΙΔΗ: Θεσσαλονίκη*, (1993), 9–14.
- Τσολάκης, Χρ. Α., *ΝΕΟΕΛΛΗΝΙΚΗ ΓΡΑΜΜΑΤΙΚΗ. ΕΚΔΟΣΕΙΣ ΚΩΔΙΚΑΣ, Θεσσαλονίκη*: (1988), 28–35.
- Dr. Wendt, H.F., *Langenscheidts Praktisches Lehrbuch Neugriechisch. Ein Standardwerk für Anfänger*, Langenscheidt: Berlin, München, Wien, Zürich, New York (1993), 20–35.

Hyphenation Patterns for Ancient and Modern Greek

Dimitrios Filippou

Kato Gatzea

GR-385 00 Volos

Greece

dfilipp@hotmail.com

Abstract

Several files with Greek hyphenation patterns for \TeX can be found on CTAN. However, most of these patterns are for use with Modern Greek texts only. Some of these patterns contain mistakes or are incomplete. Other patterns are suitable only for some now-outdated “Greek \TeX ” packages. In 2000, after having examined the patterns that existed already, the author made new sets of hyphenation patterns for typesetting Ancient and Modern Greek texts with the `greek` option of the `babel` package or with Dryllerakis’ `Greek \TeX` package. Lately, these patterns have found their way even into the `ibycus` package, which can be used with the *Thesaurus Linguae Graecae*, and into Ω with the `antomega` package.

The new hyphenation patterns, while not exhaustive, do respect the grammatical and phonetic rules of three distinct Greek writing systems. In general, all Greek words are hyphenated after a vowel and before a consonant. However, for typesetting Ancient Greek texts, the hyphenation patterns follow the rules established in 1939 by the Academy of Athens, which allow for breaking up compound words between the last consonant of the first constituent word and the first letter of the second constituent word, provided that the first constituent word has not been changed by elision. For typesetting polytonic (multi-accent) Modern Greek texts, the hyphenation rules distinguish between the nasal and the non-nasal double consonants $\mu\pi$, $\nu\tau$, and $\gamma\chi$. In accordance with the latest Greek grammar rules, in monotonic (uni-accent) Modern Greek texts, these double consonants are not split.

1 Introduction

Before 2000, one could find on CTAN four different files with hyphenation patterns for Modern Greek only, namely

- `rgrhyph.tex` by Yannis Haralambous [1],
- `grkhyphen.tex` by Kostis Dryllerakis [2],
- `gehyphen.tex` by Yiannis Moschovakis [3], and
- `grhyph.tex` by Claudio Beccari [4].

The first two hyphenation-pattern files [1, 2] are almost identical. The only difference is that the patterns by Dryllerakis contain an `\endinput` command several lines before the end-of-file. (Probably, Dryllerakis cut down Haralambous’ patterns to reduce memory usage, at a time when memory space was still rather limited.) The patterns by Moschovakis [3] are not only limited to Modern Greek, but they have been “frozen” based on an obsolete mixed US-Greek codepage for DOS and an equally obsolete \LaTeX 2.09. The end result is that some words

containing vowels with combined diacritical marks (e.g., $\epsilon\tilde{\iota}\delta\omicron\varsigma$, $\theta\epsilon\tilde{\omega}$, etc.) are not hyphenated at all.

Haralambous’ patterns [1] do not provide for the correct hyphenation of combinations of three or more consonants. In addition, they do not allow for the hyphenation of the nasal consonant combinations $\mu\pi$ (*mb*), $\nu\tau$ (*nd*) and $\gamma\chi$ (*ng*), which must be split in polytonic Modern Greek. Haralambous’ patterns erroneously split the combination $\tau\mu$ and prohibit the hyphenation of all final two-letter combinations for no apparent reason.

Beccari’s patterns [4], which are commonly used with the `greek` option of `babel`, contain a number of mistakes and are also incomplete. For example, the word $\pi\upsilon\chi\nu\acute{o}\tau\eta\tau\alpha$ is hyphenated as $\pi\upsilon\chi\text{-}\nu\acute{o}\text{-}\tau\eta\text{-}\tau\alpha$. According to some rules outlined further in this text, that word should have been hyphenated as $\pi\upsilon\text{-}\chi\nu\acute{o}\text{-}\tau\eta\text{-}\tau\alpha$. Similar bad hyphenations include $\iota\sigma\text{-}\theta\mu\acute{o}\varsigma$ (it should be $\iota\text{-}\sigma\theta\mu\acute{o}\varsigma$), $\text{\textbackslash}\lambda\chi\text{-}\mu\eta\text{-}\nu\eta$ (it should be $\text{\textbackslash}\lambda\text{-}\chi\mu\eta\text{-}\nu\eta$), etc. Beccari’s patterns also allow for separation of the consonant combinations $\delta\mu$, $\delta\nu$ and $\tau\lambda$. These

combinations should not be split, because one can find some Ancient Greek words that start with such combinations (δμῶς, δνοφερός, τλημωσύνη).

In 2000, while typesetting a large volume in polytonic Modern Greek, the author of the present article noticed the mishaps in Beccari’s hyphenation patterns and the inadequacy of all other Greek hyphenation patterns. He noticed also that hyphenation patterns for Ancient Greek, although they had been discussed by Haralambous back in 1992 [5], were not available at all in the public domain. That was the incentive for the author to revise the existing hyphenation patterns for Modern Greek and to provide in the public domain a set of hyphenation patterns for Ancient Greek.

The author has already presented these patterns in the newsletter of the Greek T_EX Friends [6, 7], but this communication is the first (and long overdue) presentation of the patterns to the global T_EX community. The patterns were created for the 1988 *de facto* Levy Greek encoding [8], which later became the Local Greek (LGR) encoding.

2 Creation of Patterns

One way to produce hyphenation patterns is by using PATGEN [9]. PATGEN scans a given database with hyphenated words and prepares a set of hyphenation patterns based on observations the programme has made. Another way of using PATGEN is modular [10]: first one creates a limited set of hyphenated words, then runs PATGEN on these words, checks the produced hyphenation patterns and expands the list of hyphenated words with those words that were badly hyphenated. The whole cycle *create word list–run PATGEN–check bad hyphenations–expand word list* is repeated until an acceptable set of hyphenation patterns is produced. To the author’s knowledge, an electronic dictionary with hyphenated words does not exist for Greek. Given the excessive morphology of the Greek words, even the modular use of PATGEN would be a daunting task. A less time-consuming effort is the translation of the simple grammatical rules for the hyphenation of Greek into patterns for the T_EX machine as it has already been done [1, 3, 4]. This is the solution chosen also by the author of the present article.

Each language has its rules and exceptions that must be duly respected. It is not rare for one language to have different hyphenation rules for different dialects, or to have different hyphenation rules for texts written in different eras. The best-known example is English, where some words are hyphenated differently depending on the continent (e.g.,

pre-face in British English and *pref-ace* in American English).

In the case of Greek, one has to distinguish—grossly—between three “dialects” that demand separate sets of hyphenation patterns:

1. Ancient Greek and old-style literate Modern Greek (*katharevousa*),
2. polytonic Modern Greek, and
3. monotonic Modern Greek.

Ancient Greek is considered essentially every text that has been written in Greek from Homeric times (8th century B.C.) to about the end of the Byzantine Empire (15th century A.D.). *Katharevousa* (literally, *the purifying*) is a formal written language (almost never spoken) conceived by Greek scholars in the period of the Enlightenment as a way to purify Modern Greek from foreign influences. It was used in Greek literature of the 19th and early 20th century, and by the Greek state from its creation in 1827 until 1976. It is still used by the Greek Orthodox Church.

Polytonic and monotonic Modern Greek are essentially the same language. The only difference is that polytonic (literally, *multi-accent*) Modern Greek uses all accents, breathings and diacritics of Ancient Greek and *katharevousa*, while monotonic (literally, *uni-accent*) Modern Greek, which was adopted officially in Greece in 1982, has just one accent mark (much to the dismay of some classicists).

The hyphenation rules for Ancient Greek and *katharevousa* have special provisos for compound words [11]. The hyphenation rules for polytonic Modern Greek make a distinction between nasal μπ, ντ and γχ (pronounced as *mb*, *nd* and *ng* respectively) and non-nasal μπ, ντ and γχ (pronounced as *b*, *d* and *g*) [12]. The hyphenation rules for monotonic Modern Greek do not distinguish between nasal and non-nasal μπ, ντ and γχ, nor do they make any special demand for compound words [13].

2.1 Patterns for Modern Greek

2.1.1 Monotonic Texts

The grammatical rules for the hyphenation of monotonic Modern Greek [13] and the corresponding hyphenation patterns are the following:

1. *One consonant between two vowels always remains together with the second vowel.* This rule can be seen slightly differently: a word is hyphenated after each vowel, for example, τη-λε-ό-ρα-ση. With the Levy character encoding [8], the corresponding hyphenation patterns are: a1 e1 h1 i1 o1 u1 w1.

2. *Double vowels* (diphthongs in Ancient Greek) that are pronounced as one are not hyphenated. Hence the double vowels α , $\alpha\acute{\iota}$, $\alpha\upsilon$, etc. should not be split apart. The corresponding hyphenation patterns are: $\mathbf{a2i}$ $\mathbf{a2'i}$ $\mathbf{a2u}$... $\mathbf{u2i}$ $\mathbf{u2'i}$. However, when the first vowel is accented, the two vowels are to be pronounced separately and they can be hyphenated. Hence, we include some exceptions: $\mathbf{'a3u}$ $\mathbf{'e3u}$ $\mathbf{'o3u}$ $\mathbf{'u3i}$.

3. *Semi-vowels are not hyphenated.* Vowels, simple and double, that are usually pronounced as i are sometimes semi-vowels, i.e., they are not pronounced totally separately from the preceding or following vowel. Hence some vowel combinations involving semi-vowel sounds (j) should not be split apart. The most common semi-vowel combinations are: aj ($\nu\epsilon\text{-}\rho\acute{\alpha}\text{-}\delta\alpha$), ej ($\zeta\epsilon\acute{\iota}\text{-}\mu\pi\acute{\epsilon}\text{-}\kappa\eta\varsigma$) and oj ($\chi\omicron\text{-}\rho\acute{\omicron}\text{-}\delta\omicron$) when they are accented on the first vowel or when they are not accented at all, and the combinations ja ($\delta\iota\alpha\text{-}\beta\acute{\alpha}\zeta\omega$), je ($\epsilon\text{-}\lambda\iota\acute{\epsilon}\varsigma$) and jo ($\text{Μ}\alpha\text{-}\rho\acute{\iota}\acute{\omega}$) when they are accented on the second vowel or when they are not accented at all. The resulting hyphenation patterns are: $\mathbf{a2h}$ $\mathbf{a2"i}$... $\mathbf{i2a}$ $\mathbf{i2'a}$... $\mathbf{u2w}$ $\mathbf{u2'w}$. Some notable exceptions are: $\mathbf{'a3h}$... $\mathbf{'u3w}$.

It is worth noting that there is an inherent difficulty in distinguishing between vowels and semi-vowels. Sometimes, two words are written the same, but they are pronounced with or without a semi-vowel, thus completely changing their meaning, e.g., $\delta\acute{\omicron}\text{-}\lambda\iota\alpha$ (the adjective *devious* in feminine singular) and $\delta\acute{\omicron}\text{-}\lambda\iota\text{-}\alpha$ (the adverb *deviously*). Distinguishing between a semi-vowel and a true vowel is very difficult and requires textual analysis [14]. For the purpose of \TeX , all such suspicious semi-vowel combinations are treated as semi-vowels. The end result is that the word $\alpha\pi\acute{\omicron}\eta\chi\omicron\varsigma$ will be hyphenated as $\alpha\text{-}\pi\acute{\omicron}\eta\text{-}\chi\omicron\varsigma$. But it is better seeing some words hyphenated with one less syllable, than seeing extra syllables in other words, e.g., $\beta\acute{\omicron}\eta\text{-}\theta\alpha$ $\text{Π}\alpha\text{-}\nu\alpha\text{-}\gamma\iota\text{-}\acute{\alpha}$! (Apparently, Liang took the same approach, disallowing some valid hyphenations for the sake of forbidding definitely invalid ones [15].)

4. *Single or double consonants at the end or the beginning of a word do not constitute separate syllables.* The corresponding patterns are $\mathbf{4b}$. $\mathbf{4g}$ $\mathbf{4y}$., $\mathbf{.b4}$ $\mathbf{.g4}$... $\mathbf{.y4}$. To these patterns, one must add some other ones for the case of elision: $\mathbf{4b''}$ $\mathbf{4g''}$... $\mathbf{4y''}$.

5. *Double consonants are hyphenated.* The patterns for this rule are: $\mathbf{4b1b}$ $\mathbf{4g1g}$... $\mathbf{4q1q}$ $\mathbf{4y1y}$.

6. *Consonant combinations that cannot be found at the beginning of Greek words must be split after the first consonant* The patterns are: $\mathbf{4b1z}$ $\mathbf{4b1j}$... $\mathbf{4y1f}$ $\mathbf{4y1q}$. No distinction is made between nasal and non-nasal $\mu\pi$ (mb/b), $\nu\tau$ (nd/d) and $\gamma\chi$ (ng/g); these consonant combinations are not to be split. However, some other patterns are inserted to deal with some thorny combinations of three or more consonants:

$\mathbf{4r5g2m}$	$\acute{\epsilon}\rho\text{-}\gamma\mu\alpha$ (Anc. Gr.)
$\mathbf{tz2m}$	$\mu\acute{\alpha}\nu\alpha\text{-}\tau\zeta\mu\epsilon\nu\tau$
$\mathbf{4r5j2m}$	$\pi\omicron\rho\text{-}\theta\mu\acute{\omicron}\varsigma$
...	
$\mathbf{4m5y2t}$	$\lambda\acute{\alpha}\mu\text{-}\psi\tau\epsilon$
$\mathbf{4g1kt}$	$\epsilon\lambda\epsilon\gamma\text{-}\kappa\tau\eta\varsigma$
$\mathbf{4n1tz}$	$\nu\epsilon\rho\alpha\nu\text{-}\tau\zeta\acute{\iota}\acute{\alpha}$
$\mathbf{4n1ts}$	$\beta\iota\omicron\lambda\omicron\nu\text{-}\tau\sigma\acute{\epsilon}\lambda\omicron$

More patterns could have been inserted here to deal with non-Greek proper names with sequences of three or more consonants transliterated into Greek. For example, the pattern $\mathbf{4r5l2s}$ could have been added to hyphenate the transliterated name *Carlson* as $\text{Κ}\acute{\alpha}\rho\text{-}\lambda\omicron\sigma\omicron\nu$ and not as $\text{Κ}\acute{\alpha}\rho\lambda\text{-}\sigma\omicron\nu$ (the latter is not allowed according to Greek grammar rules). However, the number of such words is infinite and the effort most likely worthless.

7. *Two or more consonants at the end of a word do not constitute separate syllables.* Such endings are mostly found in Ancient Greek words, or words of non-Greek origin which have become part of the Modern Greek vocabulary: $\mathbf{4k1}$. ($\pi\iota\text{-}\nu\acute{\alpha}\chi\lambda$) ... $\mathbf{4nc}$. ($\acute{\epsilon}\lambda\text{-}\mu\iota\nu\varsigma$, Anc. Gr.) ... Such words can be found easily in reverse dictionaries of Modern Greek [16].

8. *Combinations of double consonants are separated.* These are some rare combinations of non-nasal $\mu\pi$ with $\nu\tau$ and/or $\gamma\chi$ in words of non-Greek origin which are now part of the Modern Greek vocabulary, e.g., $\mathbf{4mp1nt}$ ($\rho\omicron\mu\pi\text{-}\nu\tau\epsilon\sigma\acute{\alpha}\mu\pi\rho$ = *robe-de-chambre*).

2.1.2 Polytonic Texts

The hyphenation rules that apply to monotonic Modern Greek texts apply also to polytonic Modern Greek texts. Of course, the patterns for polytonic Modern Greek had to be expanded to include all possible combinations of vowel and diacritic (breathing, accent and/or iota subscript).

As mentioned above, polytonic Modern Greek has another notable difference in hyphenation: The nasal $\mu\pi$, $\nu\tau$ and $\gamma\chi$, which are pronounced as mb , nd and ng respectively, are to be separated. On the

contrary the non-nasal $\mu\pi$, $\nu\tau$ and $\gamma\chi$, which are pronounced as b , d and g , must not be separated. In general, $\mu\pi$, $\nu\tau$ and $\gamma\chi$ are nasal, thus the patterns: $4m1p$, $4n1t$, and $4g1k$. These consonant combinations are non-nasal when they follow another consonant: $\tilde{\alpha}\lambda\text{-}\mu\pi\omicron\upsilon\text{-}\rho\omicron$, $\sigma\epsilon\beta\text{-}\nu\tau\acute{\alpha}\varsigma$, $\acute{\alpha}\rho\text{-}\gamma\chi\acute{\omicron}$, etc., or in words of non-Greek origin: $\Upsilon\text{-}\mu\pi\rho\acute{\alpha}\eta\mu$, $\mu\pi\iota\text{-}\nu\tau\acute{\epsilon}\varsigma$, etc.

For the creation of hyphenation patterns, the non-nasals $\mu\pi$, $\nu\tau$ and $\gamma\chi$ can be treated in the same way Haralambous treated Ancient Greek compound words [5]. Hence, with the help of Andriotis’ etymological dictionary [17], a list of *exceptions* was built such as:

.giou5g2k	Γιου-γχοσλάβος
5g2krant.	Βόλγχο-γκραντ
...	
.qa5n2to	χα-ντούμης
.qa5n2tr	χα-ντρῶν
.q’a5n2tr	χά-ντρα

The list of all these exceptions is quite lengthy and covers five printed pages of *Eutypon* [6].

2.2 Patterns for Ancient Greek

The grammatical rules for hyphenation of Ancient Greek are mostly the same as those for polytonic Modern Greek. Apparently, the Ancient Greeks hyphenated following the simple rule that a single consonant between two vowels in one word belongs with the second vowel: $\sigma\omicron\text{-}\phi\acute{\iota}\text{-}\zeta\omega$, $\kappa\alpha\text{-}\theta\acute{\alpha}\text{-}\pi\epsilon\rho$. The Ancient Greeks also considered non-accented words as being part of the following word [18]. For example, the Ancients would hyphenate $\acute{\epsilon}\chi\ \tau\acute{\omicron}\upsilon\tau\omicron\upsilon$ as $\acute{\epsilon}\text{-}\chi\tau\acute{\omicron}\upsilon\text{-}\tau\omicron\upsilon$. Nonetheless, rules introduced by later scholars do not allow for such extravagant hyphenations.

A very tricky rule introduced by modern scholars states that “[Ancient Greek] compound words divide at the point of union” [18]. This rule has been extended to *katharevousa* and some typographers are still using it for polytonic Modern Greek (most likely mistakenly). That rule also appears in two variations. In one variation, which has been adopted by *The Chicago Manual of Style* [19], compound words are divided into their original parts irrespective of whether those original parts have been modified or not. Therefore, one should hyphenate $\sigma\tau\rho\alpha\tau\text{-}\eta\eta\acute{\omicron}\varsigma$ ($\sigma\tau\rho\alpha\tau\acute{\omicron}\nu + \tilde{\alpha}\gamma\omega$), $\Delta\acute{\iota}\omicron\sigma\text{-}\kappa\omicron\upsilon\rho\omicron\varsigma$ ($\Delta\acute{\iota}\omicron\varsigma + \kappa\omicron\upsilon\acute{\rho}\omicron\varsigma$), etc. This is the rule followed by Haralambous for the creation of Ancient Greek hyphenation patterns for the commercial package ScholarTeX [5]. In another variation, adopted by some 19th-century scholars [20] and the Academy of Athens [21], compound words are divided into their original constituent words *only* when the first word has not lost its last vowels by elision. According to that rule varia-

tion, the word $\sigma\tau\rho\alpha\tau\eta\eta\acute{\omicron}\varsigma$ should be hyphenated as $\sigma\tau\rho\alpha\text{-}\tau\eta\eta\acute{\omicron}\varsigma$, because the first word ($\sigma\tau\rho\alpha\tau\acute{\omicron}\nu$) has lost its final $\omicron\nu$.

For the creation of hyphenation patterns for Ancient Greek, the author chose to follow the rule adopted by the Academy of Athens, because this rule has also been adopted in the manuals used in the Greek high schools and lycées [11]. Thus, with the help of two widely-used dictionaries [22, 23], a list of *exceptions* for compound words was incorporated into the list of patterns for Ancient Greek:

>adi’e2x1	ἀδιέξ-οδος
>adie2x1	ἀδιεξ-όδου
>adu2s1’w	ἄδυσ-ώπητος
>adu2s1w	ἄδυσ-ωπήτου
...	
i2s1qili’akic.	δισ-χιλιάκις, etc.
i2s1muri’akic.	δισ-μυριάκις, etc.

This list is quite extensive; it includes 1555 patterns and covers twenty-eight printed pages [24].

It is worth mentioning here that special care has been taken not to confuse Ancient and Modern Greek exceptions for the division of consonants. For example, there are no Ancient Greek words that start with the Modern Greek double consonants $\mu\pi$, $\nu\tau$, $\gamma\chi$, $\tau\zeta$ and $\tau\sigma$. Therefore, all these combinations are divided in Ancient Greek texts, with no exception. Also, combinations of stopped consonants (π , β , φ / τ , δ , θ / χ , γ , χ) and the nasals μ or ν are not divided [20].

3 Putting the Patterns to Work

The patterns have been archived in three files, which have already found their way onto CTAN [24]:

- GRMhyph?.tex for monotonic Modern Greek,
- GRPhyph?.tex for polytonic Modern Greek, and
- GRAhyph?.tex for Ancient Greek.

(The ? is a number indicating the current version.)

The first two patterns for Modern Greek were tested on a sample text created by the author, after building a new **plain** format [6]. (Incidentally, **plain** is just an extension of Knuth’s **plain** format for the use of several languages with **babel**.) The result showed considerable improvement in comparison to hyphenation results obtained by earlier set of patterns (Table 1). With another **plain** format, the hyphenation patterns for Ancient Greek were tested on five classic texts in their original:

- Herodotus, *The Histories* A, I–III;
- Xenophon, *Anabasis* A, 1.IV.11–13;
- Plutarch, *Lives* Themistocles, II.1–5;
- Strabo, *Geography*, 7.1.1–5; and

Patterns	Mistakes (%)	Misses (%)
<code>rgrhyph.tex</code> [1]	25	13
<code>grhyph.tex</code> [4]	3	16
<code>GRPhyph.tex</code> (this work)	–	3

Table 1: Results from hyphenation tests with three different sets (files) of hyphenation patterns available in the public domain. Mistakes represent erroneous hyphenations. Misses represent missed hyphenation points.

- Lysias, *Defence against a Charge for Taking Bribes*.

Surprisingly, \TeX correctly hyphenated *all* Ancient Greek words found in these texts, which cover about seven printed pages [7, 24].

The author, however, does not believe that his patterns are error-free. The Ancient Greek adjective $\pi\rho\sigma\kappa\omicron\pi\eta$ has two different etymologies and meanings: “looking out for” hyphenated as $\pi\rho\sigma\kappa\omicron\pi\eta$ ($\pi\rho\delta + \sigma\kappa\omicron\pi\acute{\epsilon}\omega$), or “an offence” hyphenated as $\pi\rho\sigma\kappa\omicron\pi\eta$ ($\pi\rho\delta\varsigma + \chi\acute{o}\pi\omicron\varsigma$). Unfortunately, \TeX does not do textual analysis and will not understand the difference. Syllables with vowel synizesis may be erroneously split apart, e.g., $\chi\rho\upsilon\sigma\acute{\epsilon}\omega$ instead of $\chi\rho\upsilon\sigma\acute{\epsilon}\omega$. Again, \TeX does not do textual analysis and it is impossible for a typesetting system to capture such small details. Finally, the use of the same patterns for typesetting a mixed Ancient and Modern Greek text will bring a few surprises. For the purpose of \TeX , Ancient and Modern Greek are better treated as two different `\languages`.

4 Creation of Patterns for `ibycus` and Ω

The patterns created by the author have already been picked up by other people who are working on other packages or systems that use different font encodings. Using a Perl script, Apostolos Syropoulos adapted the hyphenation patterns for monotonic Modern Greek and Ancient Greek for usage with Ω [25]. Using another Perl script, Peter Heslin [26] adapted the hyphenation patterns for Ancient Greek for the `ibycus` package, which can be used for typesetting texts obtained from the *Thesaurus Linguae Graecae*.

5 Conclusions

The hyphenation patterns created by the author for Ancient and Modern Greek are indeed superior to those previously found on CTAN. Nonetheless, the patterns are presently under revision to eliminate a few minor mistakes. The author anticipates that

the improved patterns will be released in CTAN very soon—probably before the TUG 2004 conference. Hopefully, these patterns will shortly after migrate into `ibycus` and Ω , and they will become the default Greek hyphenation patterns in whatever system/package becomes the successor of \TeX .

6 Acknowledgements

The author would like to thank Dr Spyros Konstantos, former lecturer of Classics at McGill University, Montreal, Canada, for giving him access to his collection of Greek dictionaries. The author would like also to express his gratitude to Karl Berry (president, \TeX Users Group, USA), Baden Hughes (University of Melbourne, Australia) and Steve Peter (\TeX Users Group/Beech Stave Press, USA) for meticulously reviewing and editing this paper.

References

- [1] Y. Haralambous. `rgrhyph.tex`, vers. 1.1. CTAN: `fonts/greek/yannis`, 1990.
- [2] K. J. Dryllerakis. `grkhyphen.tex`, GREEK \TeX , vers. 4.0 α . `ftp://laotzu.doc.ic.ac.uk/pub/tex`, 1994.
- [3] Y. Moschovakis. `gehyphen.tex`. CTAN: `fonts/greek/greectex`, 1994.
- [4] C. Beccari. `grhyph.tex`. CTAN: `fonts/greek/cb`, 1997.
- [5] Y. Haralambous. Hyphenation patterns for Ancient Greek and Latin. *TUGboat*, 13:459–467, 1992.
- [6] D. Filippou. Beltiomenoi kodikes syllabismou polytonikon kai monotonikon keimenon gia to \TeX kai to \LaTeX (Improved hyphenation patterns for polytonic and monotonic Modern Greek texts typeset by \TeX and \LaTeX). *Eutypon*, (4):1–16, 2000. In Greek.
- [7] D. Filippou. Kodikes sullabismou gia ten stoicheiothesia archaion hellenikon keimenon me to \TeX kai to \LaTeX (Hyphenation patterns for typesetting Ancient Greek texts by \TeX and \LaTeX). *Eutypon*, (5):7–15, 2000.
- [8] S. Levy. Using Greek fonts with \TeX . *TUGboat*, 9:20–24, 1988.
- [9] F. M. Liang, P. Breitenlohner, and K. Berry. PATGEN—PATtern GENERation Program, vers. 2.3. CTAN: `systems/knuth/unsupported/texware/patgen.web`, 1996.
- [10] P. Sojka. Hyphenation on demand. *TUGboat*, 20:241–247, 1999.
- [11] M. Oikonomou. *Grammatike tes Archaiais Hellenikes (Grammar of Ancient Greek)*.

- Organismos Ekdoseos Didaktikon Biblion, Athena, 5th edition, 1987. In Greek.
- [12] M. Triantafyllides. *Mikre Neoellenike Grammatike (Abridged Modern Greek Grammar)*. Aristoteleion Panepistemion Thessalonikes, Hidryma Manole Triantafyllide, Thessalonike, 2nd edition, 1975. In Greek.
- [13] *Neoellenike Grammatike (Modern Greek Grammar)*. Organismos Ekdoseos Didaktikon Biblion, Athena, 8th edition, 1985. In Greek.
- [14] G. Orphanos, Ch. Tsalides, and A. Iordanidou. Hoi hellenoglossoi hypologistes emathan na syllabizoun? (Have the Greek-speaking computers learned how to hyphenate?). In *Praktika tes 23es Synanteses Ergasias tou Tomea Glossologias*, pages 1–8, Thessalonike, 2002. Aristoteleio Panepistemio Thessalonikes. In Greek.
- [15] F. M. Liang. *Hy-phen-a-tion by Com-put-er*. PhD thesis, Department of Computer Science, Stanford University, 1983.
- [16] G. Kourmoulakes. *Antistrophon Lexikon tes Neas Hellenikes (Reverse Dictionary of Modern Greek)*. Athenai, 1967. In Greek.
- [17] N. P. Andriotes. *Etymologiko Lexiko tes Koines Neoellenikes (Etymological Dictionary of Common Modern Greek)*. Aristoteleion Panepistemion Thessalonikes, Hidryma Manole Triantafyllide, Thessalonike, 1967. In Greek.
- [18] H. W. Smyth. *Greek Grammar*. Harvard University Press, Cambridge, Massachusetts, USA, 1956. Revised by G. M. Messing.
- [19] *The Chicago Manual of Style*. The University of Chicago Press, Chicago and London, 1982.
- [20] W. W. Goodwin. *A Greek Grammar*. Ginn and Company, Boston, 1892.
- [21] I. Kallitsounakes. Orthographikon diagramma tes Akademias Athenon (The Academy of Athens Dictation Diagramm). *Praktika tes Akademias Athenon*, 14, 1939. In Greek.
- [22] H. G. Liddell and R. Scott. *A Greek-English Lexicon*. Oxford University Press, Clarendon Press, Oxford, 1968.
- [23] I. Dr. Stamatakos. *Lexikon tes Archaiais Hellenikes Glosses (Dictionary of the Ancient Greek Language)*. Ekdotikos Organismos “O Phoinix” EPE, Athenai, 1972.
- [24] D. Filippou. CTAN: `fonts/greek/package-babel/hyphenation/filippou`, 2000.
- [25] A. Kryukov. ANTOMEGA language support package, vers. 0.7. CTAN: `systems/omega/contrib/antomega/`, 2003.
- [26] P. Heslin. CTAN: `fonts/greek/package-babel/ibycus-babel/`, 2003.

The mayan Package and Fonts

Pablo Rosell-González

Facultad de Ciencias

Universidad Nacional Autónoma de México

pablo@ciencias.unam.mx

Abstract

The Mayans developed a complete positional numeral system. This fact facilitates calculation, particularly adding and subtracting. Their number system was a vigesimal (base 20) system, consisting of three different symbols: a dot which represents ‘one’, a horizontal line which represents ‘five’, and a cacao seed or stylized sea shell representing ‘zero’. Any number between 0 and 19 can be represented with these symbols, for example, 12 is represented as two horizontal lines, one above the other, and two dots on top of them.

In this paper I will present the `mayan.sty` package for typesetting Mayan numbers, as well as the `mayan*.mf` METAFONT fonts and their PostScript Type 1 variants.

The `mayan` package allows the user to generate Mayan numbers by simply invoking, e.g., `\mayan{num}` which in any display math environment will typeset the number in a vertical stack using large symbols, while in paragraph mode will separate each level with a diagonal using adequate size symbols depending on the text font size.

Optional arguments are available for boxing the numbers when stacked vertically, writing the numbers in “vertical style” inside paragraphs, and writing the Arabic values beside each position.

1 Introduction

The ancient Mayan civilization developed sophisticated knowledge in sciences, particularly astronomy and mathematics. They adopted the numeral system from the Olmecs, a vigesimal (base 20) positional system, which includes a special symbol— for some a cacao seed, for others a stylized sea shell— representing zero, primarily used for avoiding ambiguity when a position has no value. (Compare with the Babylonian numeral system, see [4].)

This *discovery* (or *invention*) is of transcendental importance, because it allowed them to make complex calculations needed for the development of very advanced astronomical issues. They developed very precise solar and lunar calendars, and were able to predict the planetary positions, particularly the position of Venus, which was very important in their religion.

One of the four surviving Mayan codices, the *Dresden Codex*, is of special interest. It contains astronomical calculations of great accuracy concerning eclipse predictions and the position of Venus. Figure 1 shows fragments from page 24 of the Dresden Codex.

It can be seen that the ‘zero’ symbol is represented with different decorations. An important future project is to create more fonts including different designs for ‘zero’ and an implementation that allows a pseudo-random choice of them.

Part of the motivation for creating this package and fonts is that the Mayan numeral system is part of the curriculum in México and other Latin American countries when positional systems are studied, and there are no tools for generating them.

On the other hand, when the `mayan` and `mayanstela` packages are finished, they could be used in conjunction with Apostolos Syropoulos’ `epi olmec` and similar packages and fonts for archaeological purposes.

2 The Mayan Numeral System

Given any integer $b > 1$ as a base, a *complete positional number system* consists of b symbols that represent the values $0, 1, \dots, b-1$ and any linear arrangement where the first entry stands for b^0 , the second for b^1 , and in general, entry n stands for b^{n-1} .

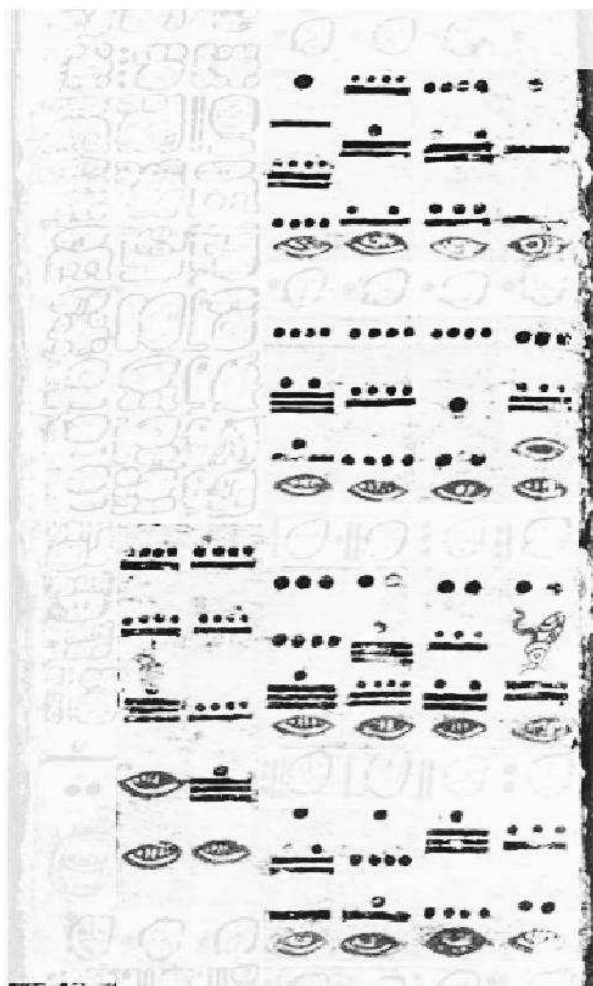


Figure 1: Fragments of page 24 of the Dresden Codex. The numbers as typically used are highlighted.

This arrangement allows us to uniquely express any number in any given base, except for the symbols used and the arrangement direction chosen. For example, if the direction of the arrangement is from left to right, the number $s_0s_1\dots s_n$ where s_i takes any value from 0 to $b - 1$ represents the number obtained by $\sum_{i=0}^n s_i \cdot b^i$.

In our decimal system, the arrangement goes from right to left, thus, any number is written as $s_ns_{n-1}\dots s_2s_1s_0$, where $s_i \in 0, \dots, 9$, and i is the entry of the arrangement. For example, 876 represents the sum of $6 \cdot 10^0$ plus $7 \cdot 10^1$ plus $8 \cdot 10^2$.

The Mayans chose 20 as the base of their numeral system, probably because their solar calendar was divided in 18 months of 20 days each, plus a five day period devoted to their *infra-world* festivities.

They used a dot to represent the unit and a horizontal bar to represent 5 units. The ‘zero’ was represented by a stylized shell and the 19 positive numbers were built according to the following rules:

1. There should be no more than four dots. (Five dots convert in a bar).
2. The dots are written above the bars.
3. The arrangement direction is vertical, from bottom to top.

From (1) and (2) we obtain, without ambiguity, the 20 required symbols:

0	☉	1	•	2	••	3	•••	4	••••
5	—	6	•—	7	••—	8	•••—	9	••••—
10	=	11	≡	12	≡•	13	≡••	14	≡•••
15	≡≡	16	≡≡•	17	≡≡••	18	≡≡•••	19	≡≡••••

Since the arrangement is vertical (rule 3), numbers greater than 19 are stacked. For example:



which stands for $14 \cdot 20^0 + 1 \cdot 20^1 + 3 \cdot 20^2 = 1234$.

It is worth noticing that each Mayan “digit” is built upon its predecessor. This allows us to add Mayan numbers without being conscious of which numbers they represent if we merely gather the respective levels and follow the rules mentioned above, along with the following:

4. Four bars at any level converts into a zero in the same level and a dot in the next.

For example:



where, for the bottom level, we only have to notice that by “passing” one dot of the left number to the right, and applying (4) we are done. The next level just consists of gathering 6 dots which converts into a bar (rule 1) and a dot above it (rule 2).

Which are the actual values of each of the above numbers? It doesn’t matter!

In other words, the symbols that represent each “digit” are actually the value of the digit in terms of bars and dots, and not an abstract glyph—in contrast to our system.

Thus, this should be considered as an alternative method for teaching the concept of *adding* at the earliest stages of education. This way the children do not have to “translate” a symbol into a set that is in bijection with the number that the symbol

represents (usually their fingers), but rather have it already written.

3 The mayan Package

The `mayan.sty` style file provides two relevant features: a base 10 to base 20 integer converter; and macros for different style layouts. When the user invokes `\mayan[opts]{num}`, the package performs two main tasks:

1. Convert num_{10} to num_{20} .
2. Depending on the context in which it is invoked, and the given optional arguments, call the required fonts and macros to display the result.

3.1 The Base 20 Converter

The classic method for converting num_{10} to base 20 is to recursively divide the number by 20 and keep the remainder. This is not useful because the number thus obtained would be written backwards. For example, let $num_{10} = 5246$, then

$$\begin{aligned} 5246 &= 20 \cdot 262 + 6 \\ 262 &= 20 \cdot 13 + 2 \\ 13 &= 20 \cdot 0 + 13 \end{aligned}$$

so $5246 = 6 \cdot 20^0 + 2 \cdot 20^1 + 13 \cdot 20^2$ but this way (without reversing the result) it would be printed as 6, 2, 13.

So, a different approach is used. We obtain p , the maximum power of 20 that is less than the given number, and divide it by 20^p . Print the result and keep the remainder. Apply this method recursively, until the remainder is less than 20. For example,

$$\begin{aligned} 5246 \div 20^2 &= 262 + 13 \\ 262 \div 20 &= 2 + 6. \end{aligned}$$

This time we obtain $5246 = 13 \cdot 20^2 + 2 \cdot 20 + 6$, which will be printed in the correct order.

The only tricky part is printing the middle positional zeros. For example, consider the number $320001 = 2 \cdot 20^4 + 1 \cdot 20^0$. The problem is solved by performing a subroutine which compares the remainder (r) with 20^{p-1} . If $20^{p-1}/r < 20$ do nothing; else print a zero and divide 20^{p-1} by 20 to obtain 20^{p-2} , and repeat until $20^{p-n}/r < 20$.

This approach is preferred because it avoids the need for creating extra macros to keep and then reverse the result. Instead we print the partial results ‘on the fly’.

3.2 Layouts

The command `\mayan` is aware of the context from which it is invoked, i.e., it will display the Mayan number horizontally and at the same font size if in-

voke inside a paragraph, or vertically and in ‘display’ size if invoked within a display math environment (see Figure 2).

As well as the default layout, `\mayan` offers the following optional arguments:

- a** Write the Arabic value to the right of each position. This option is valid only when displaying the number in vertical style. It is useful for explaining positional number systems in textbooks.
 - b** Display the number inside boxes, one for each position. Created for the same purpose as option **a**, this option is useful for emphasizing the relative values of each “digit”.
 - h** Forces horizontal style, even in display mode.
 - v** Forces vertical style, even in horizontal mode.
- Optional arguments can be given in any order.

4 The Fonts

One font family, `mayanschem`, is currently implemented for use in `mayan`, which offers simple bars and dots drawings, as well as a simple zero, mainly designed for scholarly use. These glyphs have been shown throughout the present paper.

Two more families are planned for the near future: `mayanocodex` for “handwriting” imitating the *Dresden Codex*, and `mayanstela` for “engraved” designs mimicking stelae (engraved stone and stucco), like those found in Palenque and Tikal. Both will also include diverse designs for the zero found in the codices and in stelae. The user will be able to choose specific zero designs or have the system pseudo-randomly to choose a design.

5 A Complete Example

Figure 2 shows the output for the following code:

```

1 \begin{equation}
2 \   \mayan{5246} + \mayan{3556}
3   = \mayan{8802}
4 \end{equation}
5 since \mayan{5246} equals 5246 and
6 \mayan{3556} equals 3556, as shown below:
7 \begin{equation}
8 \   \mayan[ab]{5246} + \mayan[ab]{3556}
9   = \mayan[ba]{8802}
10 \end{equation}
11 {\large then $5246+3556=8802$ is
12   represented as \mayan{8802}}
```

Lines 1–4 show a displayed equation. Lines 5–6 show text in standard layout, while lines 11–12 are `\large` text. Notice that the `\mayan` macro is invoked in the same way, no matter if math or text mode, and is aware of font size.

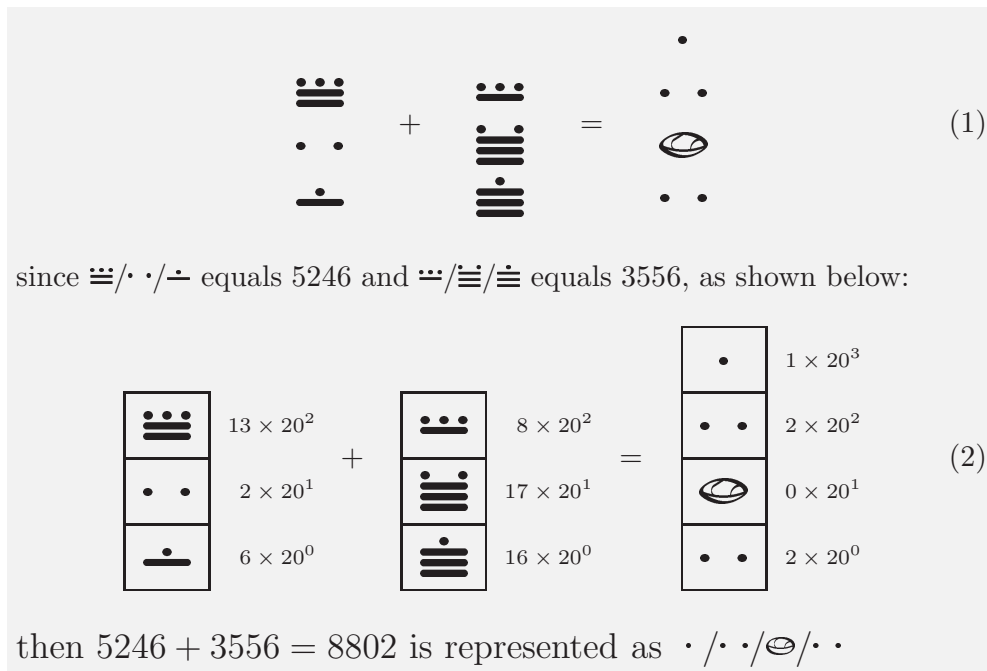


Figure 2: An example showing different layouts for mayan.

Lines 7–10 show the use of `\mayan` with optional parameters `a` and `b`.

References

[1] M. Goossens, F. Mittelbach, A. Samarin, *The L^AT_EX Companion*, Addison-Wesley, 1994.
 [2] D. E. Knuth, *The METAFONTbook*, Addison-Wesley, 2000.

[3] D. E. Knuth, *The T_EXbook*, Addison-Wesley, 2000.
 [4] O. Ore, *Number Theory and Its History*, Dover, 1988.
 [5] J. E. S. Thompson, *Códice de Dresde. Un comentario al Códice de Dresde. Libro de jeroglifos mayas*, FCE, México, 1988.

Using L^AT_EX to Typeset a Marāṭhī-English Dictionary

Manasi Athale and Rahul Athale
Research Institute for Symbolic Computation
Johannes Kepler University
Linz
Austria
manasi,athale@risc.uni-linz.ac.at

Abstract

We are using L^AT_EX to typeset an old Marāṭhī-English dictionary, dated 1857. Marāṭhī is the official language of Mahārāshtra, a western state of India. Marāṭhī (मराठी) is written using the Devanāgarī script. The printed edition of the dictionary contains approximately 1000 Royal Quarto size ($9\frac{1}{2}'' \times 12\frac{2}{3}''$) pages with around 60,000 words. The roots of the words come from many languages including Sanskrit, Arabic and Persian. Therefore the original dictionary contains at least *three* different scripts along with many esoteric punctuation marks and symbols that are not used nowadays.

We have finished typesetting 100 pages of the original dictionary. We present our experiences in typesetting this long work involving Devanāgarī and Roman script. For typesetting in Devanāgarī script we used the `devnag` package. We have not yet added the roots in other scripts but that extension can be achieved with the help of ArabT_EX. We want to publish the dictionary in electronic format, so we generated output in PDF format using pdfL^AT_EX. The bookmarks and cross-references make navigation easy. In the future it would be possible to design the old punctuation marks and symbols with the help of METAFONT.

1 Introduction

Marāṭhī is a language spoken in the Western part of India, and it is the official language of Mahārāshtra state. It is the mother tongue of more than 50 million people. It is written in the Devanāgarī script, which is also used for writing Hindi, the national language of India, and Sanskrit. The script is written from left to right. A consonant and vowel are combined together to get a syllable, in some cases consonants can be combined together to get conjuncts or ligatures. While combining the vowel and a consonant one might have to go to the left of the current character — which is a big problem for a typesetting program.

We are typesetting a Marāṭhī-English dictionary compiled by J. T. Molesworth and published in 1857. The dictionary is old so there is no problem about copyright. This will be the first Marāṭhī-English dictionary in an electronic format.

2 Devanāgarī Script

There are 34 consonants, 12 vowels, and 2 vowel-like sounds in Marāṭhī. Table 1 gives the consonants along with some common English words to

illustrate the sounds. In some cases, there is no exact equivalent English sound, and we give those with standard philological transliteration. The *h* in this table designates aspiration, and a dot under a consonant designates retroflexion. Although Hindi and Marāṭhī use the same Devanāgarī script, the consonant ष, which is used in Marāṭhī is not used in Hindi. Similarly some characters used in Sanskrit are not used in Marāṭhī. All the consonants have one inherent vowel अ (*a*), and in order to write the consonant itself without the vowel, a special “cancellation” character (.) called *virāma*, must be used. For example, स, is स् + अ, where अ is a vowel.

Table 2 lists the vowels and the two vowel-like sounds. The first two rows give the vowels and the last row gives the vowel-like sounds, called *anuswāra* and *visarga*, respectively. In general, the vowels are paired with one member short, the other long. For ऋ, make the *r* into a syllable by itself.

A vowel is added to a consonant to produce a syllable, for example all the consonants written above already have the vowel अ. Suppose we want to get the sound, *Sarah*, with a long *a*. We add the second vowel आ to स् to get सा, where we can see a bar added behind the consonant स.

क	ख	ग	घ	ङ
car	kh	go	gh	nasal
च	छ	ज	झ	ञ
chair	cch	jail	zebra	nasal
ट	ठ	ड	ढ	ण
t	th	d	dh	n
त	थ	द	ध	न
Tehran	th	dark	dh	new
प	फ	ब	भ	म
pair	fail	bat	bh	man
य	र	ल	व	
yellow	road	love	way	
श	ष	स	ह	ळ
share	ṣ	sun	happy	

Table 1: Devanāgarī consonants.

अ	आ	इ	ई	उ	ऊ
about	car	sit	seat	put	root
ऋ	ॠ	ए	ऐ	ओ	औ
under	bottle	say	by	road	loud
अं	अः				

Table 2: Devanāgarī vowels.

Now to write *sit* we add the third vowel इ to स्, and here it gets difficult, because although the *i* is pronounced after the *s*, it is written before the consonant. We get सि, where a bar is added before the character स्.

Syllables can even be formed using more than one consonant and a vowel. For example, द्वितीया, here we add द्, व् and इ. It can also be written as द्वितीया. There are many such variations when two or more consonants are combined, and some conjunct characters look nothing like their constituent parts. For example, र् or *r* is written in *four* different ways depending on the other consonant in the conjunction.

The first vowel-like sound, *anuswāra*, is the nasal consonant at the end of each of the first five consonant rows in the consonant table. For example, गंगा (Ganges), here the “ँ” on the first character is the *anuswāra* but it is pronounced as the nasal sound in the row of the next character, which is गा. The sound is like ङ. *Visarga* is more or less a very brief aspiration following the inherent vowel (and for this reason it is usually written *ḥ* in philological transcription).

3 Problems

We tried many approaches before choosing L^AT_EX. Scanning the pages was out of question as the printed quality is very poor. Also many of the consonant

conjuncts are written in a different way nowadays, so it would be difficult for the average modern reader to decipher the old dictionary. There are some web sites that have dictionaries in two scripts using Unicode. But in many cases it does not show the correct output, and it is difficult to find suitable viewers. We thank referees for mentioning an XML approach, but we did not try that. We also tried Omega, but there was hardly any information available when we started our work more than two years ago and also the setup was very difficult.

The first problem was having two scripts in the text, and typesetting it such that both scripts mesh well. Molesworth uses Marāṭhī words to explain the concepts so Devanāgarī script text appears also in the meaning. Also there are conjuncts of a poem in places to explain the usage. Many Marāṭhī words have roots in Sanskrit, Hindusthani, Arabic and Persian. Arabic, Persian, and the Urdu variant of Hindusthani are written using the Arabic script, which is the third script used in the dictionary. In Marāṭhī, a word is spoken—and also written—in a slightly different way depending on the region of the publication. Therefore in the dictionary, the most used form usually has the meaning listed for it, and all other forms have a note pointing to the most used form. This requires cross-referencing for faster use.

The dictionary has a long preface giving the details of how the words were chosen, which meanings were added, and so on. It contains different symbols and punctuation marks. Also in the meaning of some words, symbols are used to show the short form used during that period, which is obsolete now.

The printed dictionary is heavy, so carrying it everywhere is out of question. We wanted to give the user the possibility to carry the dictionary on a compact disc or computer. Therefore the next question was, which is the most user-friendly and/or popular output format?

4 Solution

In a single word: pdfL^AT_EX. We mainly used two packages to do the typesetting: `lexikon` for dictionary style, and `devnag` for Devanāgarī script. It is a two step process to typeset in Devanāgarī script. A file, usually with extension `.dn`, is processed with `devnag`, a program written in the C language, to get a `.tex` file. The preprocessing step is necessary due to the problem of vowel placement, complex conjunct characters, and so on, as mentioned in the introduction. The style file `dev` is used to get the Devanāgarī characters in the output PDF file after compiling using pdfL^AT_EX.

Once we have a `.tex` file we can get output in many formats, DVI, PS, PDF, etc. We chose PDF as there are free readers for almost all platforms and pdfL^AT_EX makes it easy to go from T_EX to PDF. The `hyperref` package solved the problem of cross-referencing and bookmarks. The user can click on a hyperlinked word to go to the form of the word that has the complete meaning, and come back to the original word with the back button in his favourite reader. In addition to the hyperlinks, bookmarks make navigation much easier; for example, bookmarks point to the first words starting with *aa*, *ab*, *ac*, etc. An additional nested level of bookmarks is chosen if there are many words starting with the character combination. For example, if there are many words starting with *ac* then we also have bookmarks for *aca*, *acc* and so on. Usually there are fewer than five pages between two bookmarks, so finding a word is not time consuming.

The preface contains characters like अण, which is not part of the modern Marāṭhī character set, but which was used as a short form a hundred years ago. To typeset this character we directly edited the `.tex` file after preprocessing to get the required result.

We have attached at the end of this article an annotated sample page from the typeset dictionary. At the top of the page the first entry is the first word on the page, then the copyright information with our name for the dictionary, शब्दविश्व, simply translated as “the world of words”, followed by the entry of the last word on the page. On the right hand side is the page number. At the bottom, the page number is given in Devanāgarī script.

5 Future Work

After completing the typesetting the whole dictionary we will add the roots of the words in Hindusthani, Arabic and Persian. Currently we denote this using [H], [A] or [P], respectively. We have tried typesetting in three scripts on some small examples and did not find any conflicts between ArabT_EX and devnag. We have not yet created new symbols but it is possible with the help of the `pstricks` package or METAFONT.

References

- [1] Devanāgarī for T_EX, <http://www.ctan.org/tex-archive/language/devanagari/velthuis/doc/devnag/manual.ps>

अकढा or अकढी *a* [P] (अ & कढणं) Unscalded or unheated—milk, oil, **ghee**. 2 Unheated to the degree of fusion—metals.

The root of the previous word is Persian, which is denoted by [P]. In the original dictionary the root is written in the original script. The word **ghee** appears in Helvetica font to stress the fact that this word is of Indian origin.

अकण *a* Devoid of कण or grit—cleaned rice. 2 wanting corn in the car—standing or thrashed crops. 3 Having no corn to eat. Ex. अकणा कण होय अधना धन होय.

अकणकडवे *n* (See आंकणकडवे) The burden or bob of a song.

The popular form of the previous word is different, which appears on the right hand side in the paranthesis. A user can click on that word to get more meanings on the current word.

अंकणी *f* A ruler. 2 (Verbal of अंकणें) Marking &c. 3 A division (as in a box), a compartment.

The previous word is a verbal of the word in the paranthesis on the right hand side. To get the meaning of the original verb user can click on that on the right hand side.

अंकणें *v c* अंकन [S] To mark, gen.; to number, stamp, dot, rule, mark with lines, figures &c.: also to rule (lines); to trace (outlines); to delineate, sketch, roughdraw, draw, describe.

The root of the previous word is Sanskrit, which is denoted by [S]. Note that Sanskrit root is given only if it is different from the word.

अंकणें *n* The sloping divisions upon a धावे or flat earthroof for the water to roll off.

अकथित *a* [S] Untold, unnarrated, unsaid.

अकथ्य & अकथनीय *a* [S] Unspeakable, ineffable, indescribable, inenarrable. Ex. परब्रह्म कीर अ°

अंकन *n* [S] Marking gen.; numbering, stamping, dotting &c.

अंकनीय *a* [S] To be marked, numbered &c. See the noun अंकन.

अकपट corruptly, अकपटी *a* [S] Free from malice or guile; forgiving, frank, ingenuous. 2 Real, true, genuine. 3 Used as *s n* Candor, openness, ingenuousness, guilelessness, absence of malice or grudge. Ex. मी तुम्हाशी अकपटाने वागतो.

अंकपट्टी *f* The ticket or label appended (to a bale of cloth &c.) showing the number and price.

अंकपाश *m* [S] In arithmetic. Permutation.

अकबरशाई *a* Of the currency established by the emperor Akbar—a rupee &c.; relating to the reign of Akbar.

Basque: A Case Study in Generalizing L^AT_EX Language Support

Jagoba Arias Pérez

Alameda Urquijo s/n

Bilbao, 48013

Spain

jtparpej@bi.ehu.es

<http://det.bi.ehu.es/~apert>

Jesús Lázaro

Alameda Urquijo s/n

Bilbao, 48013

Spain

jtplaarj@bi.ehu.es

<http://det.bi.ehu.es/~apert>

Juan M. Aguirregabiria

B^o Sarriena s/n

Leioa, 48940

Spain

wpagagj@lg.ehu.es

<http://tp.lc.ehu.es/jma.html>

Abstract

The multilingual support of L^AT_EX presents many weak points, especially when a language does not present the same overall syntactic scheme as English. Basque is one of the official languages in the Basque Country, being spoken by almost 650,000 speakers (it is also spoken in Navarre and the south of France). The origins of the Basque language are unknown. It is not related to any neighboring language, nor to other Indo-European languages (such as Latin or German). Thus, dates, references and numbering do not follow the typical English pattern. For example, the numbering of figure prefixes does not correspond to the `\figurename\thefigure` structure, but is exactly the other way round. To make matters worse, the presence of declension can turn this usually simple task into a nightmare. This article proposes an alternative structure for the basic classes, in order to support multilingual documents in a natural way, even in those cases where the languages do not follow the typical English-like overall structure.

1 Introduction

The origins of L^AT_EX are tied closely to the English language. Since those days, however, it has spread to many different languages and different alphabets. The extent of the differences among these languages is not only related to lexical issues, but to the structure of the languages themselves.

The main problem arises when the syntactic structure of the language does not follow the English patterns. In these cases the adoption of a new multilingual approach is required in order to produce documents for these languages.

Although L^AT_EX is a highly parameterizable environment, it lacks resources to alter the order of the parameters themselves. This is due to the fact that both Germanic languages (such as English and German) and Romance languages (such as French, Italian, Spanish)—and therefore the most widely spread European research languages that use the Latin alphabet—share a similar word order for numeric references. To make matters worse, the presence of declension in structure such as dates and numbers leads to a complicated generalization of procedures.

This paper describes an alternative structure for the basic classes, in order to support multilingual documents in a natural way, even in those cases where the languages do not follow the typical English-like overall structure. Specifically, the paper focuses on Basque, one of the official languages in the Basque Country, being spoken by over half a million speakers (it is also spoken in Navarre and the south of France).

The rest of the paper is organized as follows: section 2 describes the specific details of the Basque language, in section 3 a brief description of prior work is presented, section 4 describes the different approaches that can be followed to solve the problem, section 5 shows the advantages and drawbacks of the different solutions and finally, in section 6 some brief conclusions are presented.

2 Specific details of the Basque language

The origins of the Basque language are unknown. It is not related to any neighboring language, nor to other Indo-European languages (such as Latin or German). This is one of the reasons why word order and numbering schemes are different from those in English.

Dates and numbers. Basque uses declension instead of prepositions as many other languages. The main difference from other languages that use declension, such as German, is that in Basque numbers are also fully declined, even in common structures such as dates. These declensions depend not only on the case, number and gender, but on the last sound of the word. Another peculiarity of Basque is the use of a base 20 numerical system instead of the traditional decimal one.

This forces us to take into account not just the last figure of the number but the last two figures, in order to determine the correct declension for the number [3]. In the following example, two dates are represented using ISO 8601 and its translation into Basque.

2004-01-11 : 2004ko urtarrilaren 11n
2005-01-21 : 2005eko urtarrilaren 21ean

Note that although both days end in the same figure, the declension is slightly different. The same happens to the year. The extra phonemes have been added to avoid words that are difficult to pronounce. This makes automatic date generation difficult, because it must take into account all the possible cases (as base 20 is used, there may be as many as 20 different possibilities). The different number endings are shown in table 2. Note that there are only twenty

Table 1: Endings

Number	Ending (year)	Ending (day)
00	ko	-
01	eko	ean
02	ko	an
03	ko	an
04	ko	an
05	eko	ean
06	ko	an
07	ko	an
08	ko	an
09	ko	an
10	eko	ean
11	ko	n
12	ko	an
13	ko	an
14	ko	an
15	eko	ean
16	ko	an
17	ko	an
18	ko	an
19	ko	an
20	ko	an

possible terminations, and two declension classes are necessary.

Word order When numbering a certain chapter, section, etc., in English-like languages the order is always the following: first, the item class (e.g. “figure”) is named and, afterwards, the number is written. For example, we have “Figure 1.1” or “Table 2.3”. However, this is not the case in Basque. In this language, we must reverse this order: “1.1 Irudia” or “2.3 Taula”. The same applies for chapters, sections and other kind of text partitioning structures.

3 Related Work

Multilingual support for L^AT_EX is traditionally performed using the Babel package [2]. In this package, the overall structure of documents, such as books, articles, etc., is fitted to different languages by using different variables for the different strings in each language.

For example, we can take the way figure captions are numbered in these types of documents: a variable called `\figurename` contains the string corresponding to the word “figure” in the first part of the caption, while another variable, `\thefigure` contains the number assigned to that caption. When a new figure is inserted in the document, the string preceding the caption is always formed by using a

concatenation of both variables. However, this process is not performed by Babel, which would allow a general description of the language, but in the different files that describe the document format: `book.cls`, `article.cls`, etc. Thus, some of the work that should be performed by the module in charge of the multilingual support is made by the formatting part of the typesetting software.

The file `basque.1df` [1] currently provides support for Basque in Babel. In this file, the most commonly used words have been translated. However, this does not solve the problem of the different order of strings. In [1], a possible solution is proposed using a new package for the document definition: instead of using the multilingual capabilities of Babel to solve the problem, a new document formatting file is defined, where the specific corrections for the language are performed. The limitation for multilingual document generation is obvious in this scheme: the format must be redefined whenever the language of the document is changed. Besides, a new definition for every single class of document must be performed for this particular language — as we are not philologists, we do not know if the same happens in other languages.

4 Approaches to the Solution

The solution to the problem described in this paper must deal with the following issues:

- It must respect all the translations of the different strings generated automatically.
- It must respect not only the translation, but the order of words as well.
- The last problem to solve is the use of the `\selectlanguage` directive, which would allow us to change the hyphenation patterns and the automatic text generation structures dynamically in the same document. This directive is particularly useful for documents which contain the same text in different languages (e.g. user's guides, where the manual has been translated).

The main possible avenues to the solution are the following:

- **Use of specific classes for the language:** This solution implies the redefinition of every document format, in order to embed the corresponding word order alteration for automatic string generation. The main drawback of this alternative is the need for rewriting and adapting all the existing document formats.
- **Use of a specific package for the language:** A second possibility could include the definition of a new package for those languages that

require a word order alteration. This package should redefine the `\fnum@figure` and the `\fnum@table` variables (among others, which define the chapter or section name) in order to adapt them to the needs of the languages used. A macro should be used to switch between the two nodes.

- **Inclusion of order parameters in the document class definition files:** This option requires that a new input parameter is defined in the document class to define the order of the words. Basically, it is the same solution as the first one, but merging all the different files for a document class into a single (larger and more complex) file.
- **Redefinition of existing multilingual support files:** This solution implies the addition of several lines to every language support file, where the definition of the automatic strings such as the figure captions or the table captions is performed. For example, for the case of table and figure captions, the definitions for the Basque language would be the following:

```
\def\fnum@figure{\thefigure~\figurename}
\def\fnum@table{\thetable~\tablename}
```

These definitions should go into the `basque.1df` file, immediately after the definition of the terms for caption or table names. Thus, whenever a `\selectlanguage` directive is introduced in the document, the Babel package will read the definitions for the new language, which will include the definitions for every string.

5 Comparison of Solutions

We use the following criteria to compare the different solutions:

- **Extent of modification to existing files:** This criterion measures how many existing files will be altered to fix the problem and how complicated this alteration is.
- **Addition of new files:** This criterion measures how many new files are to be added to the L^AT_EX distribution for each solution.
- **The `\selectlanguage` issue:** This criterion measures how well the solution deals with possibly changing the language of the document dynamically.
- **How easily new automatically-generated strings are included:** In the future, translation of new strings may be required. Therefore, the proposed solution must provide an easy way to include these new strings.

5.1 Extent of Modification

Here is how the solutions fare with respect to the first criterion:

- **Use of specific classes for the language:** This option does not require that any file be modified, because new definitions are described in new files.
- **Use of specific package for the language:** This approach requires no modifications of existing files, since all modifications are included in a new package.
- **Inclusion of order parameters in the document class definition files:** This alternative entails the redefinition of every document class. These should admit a language parameter to determine the correct word order.
- **Redefinition of existing multilingual support files:** This choice implies that every file containing the translation and definition of the automatically-generated strings provides order information for them, and therefore, all the files in Babel should be changed.

5.2 Addition of New Files

Here's how the solutions fare with respect to adding new files:

- **Use of specific classes for the language:** This option requires all document classes to be rewritten for every language that does not follow the English-like structure.
- **Use of specific package for the language:** This approach requires one new file for every language that has not been described successfully in the Babel approach.
- **Inclusion of order parameters in the document class definition files:** This alternative entails no new files, as it is based on the modification of the existing files.
- **Redefinition of existing multilingual support files:** This choice does not need new files, as it is based on the modification of the existing files.

5.3 The `\selectlanguage` Issue

Depending on how generalization of the multilingual support is implemented, the different solutions may (or not) solve the `\selectlanguage` problem:

- **Use of specific classes for the language:** This option does not really use Babel and its macros. As part of the translation of automatic strings is performed by the file defining the format of the document class, support for

the `\selectlanguage` directive should be implemented in each document class for every language (not only for those incorrectly supported by the Babel system, but for all of them).

- **Use of specific package for the language:** This approach requires one new file for every language. Hence, a macro would be required in each package to leave things as they were *before* the package was initiated.
- **Inclusion of order parameters in the document class definition files:** This alternative cannot solve the problem, because the order specification is only made at the beginning of the document. A macro could be added to alter its value dynamically throughout the document, but it would be an artificial patch that would not fit naturally in the Babel structure.
- **Redefinition of existing multilingual support files:** This choice does solve the problem, because when a new `\selectlanguage` command is issued, the definitions for the new language are reloaded. It requires no new macro definitions to suit the Babel scheme for multilingual documents.

5.4 Inclusion of New Strings

Here's how the solutions fare with respect to the possibility of including further modifications for strings that could be necessary in the future:

- **Use of specific classes for the language:** As some of the linguistic characteristics of the document are included in the document class, this option does not provide a straightforward method for including changes for problems that may arise.
- **Use of specific package for the language:** The use of a package gives flexibility to the scheme, allowing the insertion of new macros to adapt to the peculiarities of the language. However, the range of possibilities is so wide that a very well-defined structure must be laid down in order to keep a modicum of coherence for creating a document in a different language.
- **Inclusion of order parameters in the document class definition files:** This scheme requires updating several files whenever a new string or scheme must be added.
- **Redefinition of existing multilingual support files:** As this choice uses a single file for every language, it makes updating the elements for Babel very easy.

Table 2: Solution comparison

Solution	Mod.	Cr.	Multi.	Updates
Specific class	X	✓	Dif.	Dif.
Specific pack.	X	✓	Dif.	Dif.
Parameters	✓	X	Dif.	Dif.
Redefinition	✓	X	✓	✓

6 Conclusions

This paper discusses some alternatives to solve the ordering problems that may arise in multilingual documents.

The characteristics of the different proposed solutions are summarized in table 2. Among the solutions, the most suitable would be the redefinition

of all the existing Babel files. The reason is simple: it requires the addition of two lines to approximately 45 files, and allows the update of the system in the future, as it maintains *all* the translating issues within their natural context (Babel).

References

- [1] Juan M. Aguirregabiria. Basque language definition file. <http://tp.lc.ehu.es/jma.html>, 2001.
- [2] Johannes Braams. Babel, a multilingual package for use with L^AT_EX's standard document classes. CTAN://macros/latex/required/babel/, 2001.
- [3] Euskaltzaindia. Data nola adierazi. <http://www.euskaltzaindia.net/arauak/dok/ProNor0037.htm>, 1995.

Implementation Tricks in the Hungarian Babel Module

Péter Szabó

Department of Analysis
Budapest University of Technology and Economics
Műegyetem rakpart 3–9.
Budapest H-1111
Hungary
pts+tug@math.bme.hu

Abstract

`magyar.ldf`, the Hungarian Babel module, was rewritten in the autumn of 2003 to obey most of the Hungarian typographical rules. This article describes some implementation issues, \TeX macro programming hacks, and \LaTeX typesetting trickery used in `magyar.ldf`. All features of the new `magyar.ldf` are enumerated, but only those having an interesting implementation are presented in detail. Most of the tricks shown are useful for developing other language modules.

1 The Name of the Language

Usually a Babel language module has the English name of that language. For example, the German module is called `germanb.ldf`, and not `deutsch.ldf`. The Hungarian module is an exception to this rule, because it has the name `magyar.ldf`, in which “magyar” is the Hungarian adjective meaning “Hungarian”. A similar exception is `portuges.ldf` for Portuguese. The letter “a” in word *magyar* has to be pronounced as in *blah*, and the consonant “gy” is the same as “d” in *due*.

The name of a language that a Babel language module (`.ldf` file) defines is usually specified as an argument of `\LdfInit` in the file. Thus, if `czech.ldf` is renamed to `foo.ldf`, it will have to be loaded with `\usepackage[foo]{babel}`, but to activate it, `\selectlanguage{czech}` should be used. This is not the case with `magyar.ldf`, because it detects its loaded filename using the `\CurrentOption` macro set by the `\ProcessOptions` command called from `babel.sty`. So whatever `magyar.ldf` is renamed to, that name is the one to pass to `\selectlanguage`.

The only reason why someone may want to rename an `.ldf` file is to load two different versions in the same \LaTeX run. This is possible with `magyar.ldf`, but the user should be aware that the control sequences defined by the two copies will interact with each other in an unpredictable way. Experiments have shown that it is possible to load two copies of `magyar.ldf` with different load options (this is the so-called dual load):

```
\PassOptionsToPackage{frenchspacing=yes}
{magyar.ldf}
```

```
\PassOptionsToPackage{frenchspacing=no}
{hungarian.ldf}
\usepackage[hungarian,magyar]{babel}
```

Despite the name `hungarian.ldf` above, the file `magyar.ldf` gets loaded twice, because Babel translates the language name `hungarian` to file name `magyar.ldf`, and `magyar.ldf` expects options for `\CurrentOption.ldf`, which depends on the language name passed to `\usepackage[...]{babel}`. Since the dual load feature of `magyar.ldf` is experimental, most of the load options cannot be different in the two copies. So the safest way to load two copies is to replace the occurrences of the word `magyar` in the second copy with something else.

The latest `magyar.ldf` (version 1.5) is not part of standard Babel yet, but it is available as part of `Magyar \LaTeX` (see section 4.1). Most of the typographical rules it tries to obey and problems it addresses were proposed in [1].

2 What an `.ldf` File Contains

An `.ldf` file is a Babel language module, which contains specific macros for the given language. It is loaded by `babel.sty` in the document preamble, at the time `babel.sty` itself is loaded. The macros defined in `foo.ldf` take effect only after changing the language with `\selectlanguage{foo}`. The default language is the one specified last in the `\usepackage[...]{babel}` command.

Babel itself contains the standard versions of the `.ldf` files as `tex/generic/babel/*.ldf`. In Babel 3.7 there are 41 of them; most are smaller than 10 kB. The largest files are: the old `magyar.ldf` defining the Hungarian language (25 kB), `frenchb.ldf`

defining the French language (23 kB), `spanish.ldf` (21 kB), `bulgarian.ldf` (13 kB), `ukraineb.ldf` defining the Ukrainian language (12 kB), `russianb.ldf` (12 kB) and `greek.ldf` (9 kB). The new version of `magyar.ldf` is much larger than any of these: it is 178 kB. The size implies much more functionality, including several features unique to this new `magyar.ldf`—they will be discussed later in this document. Let’s proceed first by dealing with features common in most `.ldf` files.

2.1 Selecting the Hyphenation Pattern Set

`foo.ldf` must define the control sequence `\l@foo` to be a number (`\newcount`, `\chardef`, etc.) representing the hyphenation pattern set to be used for that language. The language selection macro `\selectlanguage{foo}` calls `\language=\l@foo`, which activates the hyphenation patterns for the language `foo`.

The patterns were (presumably) defined with the `\patterns` primitive at the time `iniTeX` was called to generate the format file. The exact filename containing the `\patterns` command is specified in the file `language.dat`. If there is a line “`foot fthyphen.tex`” in `language.dat`, then `\language=\l@foot` will activate `\patterns` found in `fthyph.tex`. In that case, `foo.ldf` should contain a line `\let\l@foo\l@foot`. But this line is omitted in most actual `.ldf` files, because the Babel language name and the hyphenation pattern set name is the same (`language.dat` would contain an entry starting with `foo_` in our example). Note that the file `fthyph.tex` is read by `iniTeX`, not `LaTeX`, so the format files have to be re-generated each time `fthyph.tex` is changed.

Three different hyphenation pattern sets have been proposed for the Hungarian language (namely, `huhyp3.tex`, `huhypc.tex` and `huhypf.tex`). All of them are maintained by Gyula Mayer [4]. The most important difference among them is the way they hyphenate at subword boundaries of compound words. The document author can select any of these three by providing the appropriate load option to `magyar.ldf` (discussed later). The options work by redefining `\l@magyar` to be one of `\l@magyarf`, `\l@magyarc`, or `\csname l@magyar3\endcsname`.

There are two different correct ways to hyphenate compound words in Hungarian. `magyarf` hyphenates the most common foreign compound words of Hungarian text phonetically (e.g. *szink-ron*, meaning synchronous), while `magyarc` hyphenates them on the subword boundary (e.g. *szin-kron*). `magyar3` is the old version of the hyphenation patterns which hyphenates most composite words phoneti-

cally (even non-foreign ones), save only a few exceptions listed explicitly. However, in all the three cases, hyphenation of foreign words cannot be perfect, because *all* of them cannot be specified in `\patterns`.

`magyar.ldf` redefines `\l@magyar` depending on the `hyphenation=` load option. If a given pattern set may be missing from the user’s system, `magyar.ldf` falls back to another set with a meaningful warning message. Hyphenation is disabled not by choosing `\language0`, as Babel does, because `\language0` may contain valid patterns for a different language, but rather `\language255`, which is very likely to be unused since `LaTeX` assigns `\language` numbers from zero.

2.2 Defining Captions

`LaTeX` generates some words and phrases automatically. For example, `\tableofcontents` should emit the phrase “Table of contents” in the native language. The same applies for `\captions` of figures and tables, and also for `\chapter` titles. Thus Babel expects `foo.ldf` to define a macro called `\captions foo` containing definitions like `\def\abstractname{Absztrakt}`. These definitions are executed by `\selectlanguage` each time the language is activated. So it is possible to have an English and then a Hungarian chapter in a book numbered ‘Chapter 1’ and ‘2. fejezet’, respectively:

```
\chapter{foo} ...
\selectlanguage{magyar}\chapter{bar}
```

`magyar.ldf` has the proper definitions of Hungarian phrases. Some words contain accented letters, which are specified as commands (e.g. `\’a` for *á*) and not single 8-bit characters, so their interpretation does not depend on the active input encoding, i.e. the load option of `inputenc.sty`.

2.3 Generating Dates

`foo.ldf` should define a macro `\datefoo` to define the macro `\today`, which emits a date (specified by the `\year`, `\month`, `\day` registers) correctly for that language. The month name should be printed as a non-abbreviated word. The definition of `\today` is used by `\@date` invoked in `\maketitle` in the standard document classes.

In addition to defining `\today`, `magyar.ldf` defines the macro `\ondatmagyar` to further define `\ontoday`, which emits the date with the Hungarian equivalent of English *on*. The Hungarian language has suffixes instead of prepositions, and each suffix has several forms which must follow the vowel harmony of the word it is suffixed to. Thus “on March

15” is emitted as *március 15-én*, but “on March 16” is *március 16-án*, showing that the *-án/-én* suffix has two forms.

2.4 Minimum Hyphenation Length

TeX won’t insert an implicit hyphen into the first `\lefthyphenmin` characters of words, nor in the last `\righthyphenmin` characters. The default TeX values for these are 2 and 3, respectively, which are suitable for the English language. `foo.ldf` can override the default by defining the macro `\foohyphenmins` to be *lr*, two digits specifying the left and the right minimum, respectively.

What `magyar.ldf` does depends on its load options. The default is to follow Hungarian typography: `\def\magyarhyphenmins{22}`.

Nine of the 41 `.ldf` files in Babel 3.7 do only the customizations described to this point. 25 languages go a little beyond these, and 7 languages go much beyond. Of those 25 + 7 languages that go beyond, we will compare `frenchb.ldf` in detail to `magyar.ldf`, because French and Hungarian share some typographical rules.

2.5 Defining Special Letters

Many languages have letters that are missing from the standard OT1 encoding, and some characters are missing even from T1. These should be implemented in `.ldf` files as control sequences. It is a common practice to modify the meaning of an existing letter, for example `czech.ldf` contains `\DeclareTextCompositeCommand{\v}{OT1}{t}{...}`. However, this declaration is contained in `\AtBeginDocument`, so they are in effect even when not the Czech language is active. This should have been avoided.

The correct solution is to use the *extras* facility provided by Babel: `foo.ldf` can have a macro `\extrasfoo`, which is executed each time the language `foo` is activated; and the macro `\noextrasfoo` is executed when the active language is about to change (because of a `\selectlanguage` command or when the end-of-group is reached). It is a common practice in `\extrasfoo` to save the meaning of a macro with `\babel@save`, or a meaning of a count, dimen or skip register with `\babel@savevariable`. The saved meanings will be restored just after `\noextrasfoo` is executed. Babel provides the command `\addto` that can append tokens to the definition of an existing macro. The idiom `\addto\extrasfoo{\babel@save\bar \def\bar{foo-bar}}` is typical, which gives a new meaning to `\bar` while the language `foo` is active.

The macro to be saved for `\DeclareTextCompositeCommand{\v}{OT1}` is `\OT1\v` (with the second backslash being part of the control sequence), but assigning the new meaning would be problematic, since `DeclareTextCompositeCommand` can be used only in the preamble. Thus the correct solution would involve fiddling with undocumented L^AT_EX internals; which is probably why `czech.ldf` contains the problematic workaround using `\AtBeginDocument`.

Fortunately, the only non-English letters in the Hungarian language are accented vowels (*á, é, í, ó, ő, œ, ú, ű* and *ű*), which are all part of the T1 encoding. The letters *ő* and *ű* with the special Hungarian double-acute accent are missing from the Latin-1 encoding (ISO-8859-1), but are part of Latin-2. So authors dealing with Hungarian are encouraged to use `\usepackage[latin2]{inputenc}`.¹ `\usepackage{t1enc}` is also recommended, so TeX will be able to hyphenate words containing accented letters.

The finest Hungarian books have accents lowered a little bit. This is accomplished for the dieresis accent (¨) by calling the `\umlautlow` command (defined by Babel) in `\extrasmagyar`. No serious attempt is made to make this work for all three Hungarian accents, because the technology `\umlautlow` is based on works only for the OT1 encoding (which composes accented letters), but most Hungarian texts use the T1 encoding to allow hyphenation in words with accented letters.

The lowering of accents is possible using virtual fonts. But TeX font families come with too many variations and design sizes, so the virtual font generation would need to be automated. The macro `\lower@umlaut` in `babel.def` lowers accents by forcing their top to be 1.45ex above the baseline. The `\accent` primitive lowers its accent by `\fontdimen5\font-1ex`, so the top of the accent can be forced to 1.45ex by setting `\fontdimen5\font:=\ht0 - 0.45ex`, where `\ht0` is the height of the accent character (`\char127` in the OT1 encoding).

The lowering, in the case of *ű*, is as small as 0.43558 pt. Even this tiny displacement can make a visible difference: “*ű* < *ü*”. The lowering method could be made adaptive by rendering the glyphs involved at high resolution, measuring the number of pixels between the accent and the letter vertically, and then lowering the accent so the distance will be a prescribed constant value.

¹ The most common incorrect letters found in Hungarian texts are *ő* and *ű*: their presence is caused by software incapable of using Unicode or the Latin-2 encoding. These letters can be seen even on some huge advertisement banners on streets in Hungary. These texts were not typeset by TeX, of course!

Neither home users nor professionals use lowered accents in Hungary today, not even with books created with L^AT_EX — the original fonts with the T1 encoding are acceptable enough not to bother changing. TypoT_EX Ltd., one of the biggest Hungarian publishing houses using T_EX, developed the OM fonts in the early 1990s for use with plain T_EX. The OM fonts are a variation of CM fonts with Hungarian accented glyphs added (with lowered accents). However, it is not worth creating .fd files for the OM fonts for use with L^AT_EX, because with the same amount of work new virtual fonts could be created from the EC fonts, which would take advantage of the full T1 character set, and existing, hinted fonts in Type 1 formats (such as the CM-Super fonts).

2.6 Hyphenation of Long Double Consonants

Hyphenating long double consonants in Hungarian is a difficult typographical problem. For example, the correct way to hyphenate the double consonant *tty* = *ty* + *ty* in the word *hattyú* (“swan”) is *haty-ttyú*. (There is a similar problem in German with words containing *ck*; [5] documents more languages with more exceptions.) The long double consonants involved are: *ccs*, *ddz*, *ddzs*, *ggy*, *nny*, *ssz*, *tty* and *zzs*. T_EX’s automatic hyphenation algorithm cannot deal with such exceptions, but adding ligatures dealing with the dash inserted by the implicit hyphenation can solve the problem. The simple trick of having `\patterns{t1ty}` and `t+- → ty-2` seems to solve the problem, because it hyphenates *tty* as *ty-ty*, but it also inserts an extra *y* before the hyphen in *fut-szalad*. Normal patterns will also insert an implicit hyphen into *botcsinálta*, yielding *bot-csinálta*. The ligature program above would then incorrectly alter that to *boty-csinálta*.

So a more elaborate set of ligatures would have to be constructed, to detect the context of the hyphen and insert the *y* only into *t-ty*, yielding *ty-ty*. Or, equivalently, using `\patterns{tt1y ggy1y}` with context-sensitive ligatures changing *tt-y* to *ty-ty* and *gg-y* to *gy-gy*, etc. This solution uses up many character positions from the font, and many many extra ligatures are involved. Also, the user must know that to produce an actual *t-ty* (which almost never appears in Hungarian), `t{}-ty` must be used.

All of this can be accomplished using virtual fonts. The author has tested to see that the concept works by decompiling *aer10.vf* to *aer10.vpl* and modifying the (LIGTABLE). However, automation of

the virtual font generation is work remaining for the future.

Hyphenation of the double two-character consonants *ggy* and *ssz* is similar to *tty*. However, compound words such as *leggyakoribb* (“most frequent”) and *vasszekér* (“iron chariot”) should be hyphenated at the subword boundary without the addition of extra letters, i.e. as *leg-gyakoribb* and *vas-szekér*. Extra `\patterns` may be added, for example `\patterns{gg1y .leg1g4yakoribb.}`, to disable insertion of *y* for each important compound word. This is quite straightforward, because it does not require more ligatures (apart from the context sensitive ligature program changing *gg-y* to *gy-gy*).

One might suggest context-sensitive ligatures could be avoided if *ty* is introduced as a new single letter. But this won’t work because, step (1): ‘*t y*’ has to be converted to ‘*ty ty*’ using more than one ligature, and then step (2): further conversion to ‘*t ty*’ if there is no line break, but T_EX won’t run its hyphenation algorithm in the middle of ligature processing, between steps (1) and (2).

The current approach of *magyar.ldf* for handling long double consonants is a compromise. By default, the patterns do not hyphenate those consonants, and the character ‘ is made active (with the standard Babel command `\declare@shorthand`), so that, for example ‘*tty* emits `t\nobreak \discretionary{y-}{}{}ty \nobreak \hskip \z@skip`. The first `\nobreak` is used to enable automatic hyphenation before the ‘*tty* construct, and the last `\nobreak` plus the `\hskip` enable hyphenation after the construct. Thus the word typed as *megho’sszabbít* will be hyphenated as *meg-hosz-szab-bít*. Similar shorthands are added for the other long consonants. The compromise is that the user has to be aware that he has to insert ‘*s* manually. A Perl script named *ccs_extract.pl* was developed to collect all occurrences of double consonants in a document so the user can review them and decide about the insertion of ‘ for each.

2.7 Table and Figure Captions

The document class defines `\@makecaption`, which is responsible for typesetting a caption for tables and figures. Some .*ldf* files, including *frenchb.ldf* and *magyar.ldf* override the default behaviour. *magyar.ldf* changes the colon separating the caption heading (e.g. “1. táblázat”, or “Table 1”) from the caption text to a full stop, in keeping with Hungarian typography. Furthermore, the `longcaption=` load option controls what should happen when the caption doesn’t fit in a single line: whether it should

² (LABEL C t) (LIG C - C y) in the .pl file.

be centered, and whether there should be a line break after the caption heading.

The `tablecaptions=` and `figurecaptions=` load options control the appearance of the caption heading by redefining both `\fnum@table` and `\fnum@figure`. The default in both cases is to follow Hungarian typography, which requires the number to precede the table name.

2.8 Between the Section Title and the Section Number

The default definitions of `\@ssect` and `\@sect` separate the section number from the section title with a `\quad`. In Hungarian typography, only an `\enskip` is needed, and a dot has to be inserted after the number. The old version of `magyar.ldf` changed `\@sect` etc., but this caused conflicts with the AMS document classes and other packages, so that strategy has been given up in the new version of `magyar.ldf`. Instead, dots were moved into the `\numberline`, which adds the dot to `\tableofcontents` lines (being careful not to append the dot twice, see dot stripping code in section 2.26), and to `\@secntformat`, which adds the dot and `\enskip` to the titles. The AMS document classes do not use `\@secntformat`, so the AMS-specific `\tocsection` and `\tocappendix` commands had to be modified.

`\@numberline` also adds a dot after table and figure numbers in the `\listoftables` and `\listoffigures`, but the dot is needed there, too, anyway.

All three TOCs share a common problem related to language changes. Each time the language is changed, `Babel` emits changing commands to the three TOC files, so when they are re-read, each line comes in its appropriate language. The implementation has a flaw, however, because the `\write` to the TOC files gets executed when the page is shipped out, and the order of `\writes` on the same page follows the document structure: `\writes` in top insertions precede and bottom insertions follow those in the main text. So when a table or figure is moved to the top of the page, the writing of its TOC entry, together with the `\selectlanguage` command emitted by `Babel` is moved away, so `\selectlanguage` commands in the TOC files are reordered, which is wrong. The solution is to emit a `\selectlanguage` command for each TOC entry, so the TOC entries can be freely reordered. `magyar.ldf` implements this solution as a local fix, but it should be fixed generally in a new version of `Babel`.

2.9 Spacing Around Punctuation

It is quite easy to add extra space *after* punctuation characters with `\sfcode` (see “space factor” in chapter 13 of [3]). The L^AT_EX `\nonfrenchspacing` command (which is activated by default) assigns a space factor of 3000 to `.`, `?` and `!`, 2000 to `:`, 1500 to `;`, and 1250 to `,`.

However, adding extra space *before* punctuation needs a different approach. Both `frenchb.ldf` and `magyar.ldf` make the characters `:`, `;`, `!` and `?` active with the `Babel \initiate@active@char` interface, and insert unbreakable space in horizontal mode (`\ifhmode`) just before the punctuation character. This feature of `magyar.ldf` can be turned off using the `activespace=` load option, partly because making these four common characters active may lead to incompatibility with other packages, and partly because the extra space before punctuation is very rare in current Hungarian documents. In French typography, about one-third of a normal space is required before punctuation, and if it is not possible to add that amount with the typesetting technology, one full space should be added. However, in Hungarian, the fallback strategy is to omit the extra space.

The last action the active punctuation character should do is insert itself, but typing it verbatim into the definition will lead to an infinite loop. For example, `\catcode'?=13 \def?{\kern.1em ?}` will loop infinitely. The solution is to use `\string?` in place of the last `?`, so its catcode will be changed to 12 (other). Using `\edef` with this approach will make the macro a little bit faster, because `\string` will be executed only once, at load time.

2.10 Quoted Material

In English, text can be quoted using ‘single’ or “double” quotation marks. These can be nested into each other both ways. Hungarian provides three nesting levels: „outer »middle ’inner’ «”. Although the guillemet symbols are missing from the CM fonts with OT1 encoding, this is not a serious problem, since `Babel` provides them (using the `\ll` relation: `<<` and `>>`), and Hungarian text should be typeset with the T1 font encoding anyway, to allow hyphenation of words with accented characters.

`frenchb.ldf` provides `\LasyGuillemets` and `\CyrillicGuillemets` so the user can select the origin of the replacement guillemets. `magyar.ldf` relies on the defaults provided by `Babel` in the hope that the T1 encoding is used, so replacements are not needed. `magyar.ldf` doesn’t adjust spacing around

the quotation symbols, but provides a `\textqq` command which emits quotations with proper nesting and spacing. For example, `\textqq{outer \textqq{middle \textqq{inner}\,}\,}` gives the above three-level sample. `\textqq` does English quotations (with two alternating levels) when the Hungarian language is not active.

2.11 List Environments

The default spacing, indentation and label item generation of list environments (such as `itemize` and `description`) are incorrect for Hungarian. The `labelenums=` and `labelitems=` load options control whether labels are modified to match Hungarian traditions. Five levels of depth are provided for both `itemize` and `enumerate`. The maximum depth is hardwired to the \LaTeX definitions of these environments, so the `\ifnum\@enumdepth>3` test had to be changed to

```
\expandafter \ifx\csname \labelenum
\romannumeral\the\@enumdepth\endcsname\relax
( and similarly for \@itemdepth).
```

Although the vertical space that the standard document classes leave around lists is too large, and the indentation is also incorrect, these problems have not yet been solved in `magyar.ldf`. (`frenchb.ldf` modifies `\itemsep` and other spacing dimensions to match French typographical rules.)

2.12 Modularity Using Load Options

The user can customize `.ldfs` using Babel's language attribute facility. For example, `greek.ldf` has `\bbl@declare@ttribute{greek}{polutoniko}{.}`, and so if the user loads `greek.ldf` with

```
\usepackage{greek}{babel}
\languageattribute{greek}{polutoniko}
```

the code in the `..` is run when `\languageattribute` is called. If present, `\languageattribute` must be part of the document preamble, and the `.ldf` file must be already loaded.

The fundamental problem with language attributes is that the user can pass only declared keywords, and not arbitrary data to the `.ldf` file, and — since attributes are processed too late — they cannot be used to control which parts of the `.ldf` files should be loaded.

Thus, `magyar.ldf` follows a different approach. Options are `<key>=<value>` pairs, which can be declared any time before `magyar.ldf` is loaded. The set of keys is fixed, but values can be arbitrary. It is the responsibility of the macro belonging to the key to verify that the syntax of the value is correct. None of the other `.ldf` files provide load option support this flexible. This option-passing scheme is similar to `keyval.sty`, but `magyar.ldf` doesn't actually use `key-`

`val.sty`, because of a general design policy to avoid dependencies.

Since `.ldf` file names are the \LaTeX options to `babel.sty` in the `\usepackage[...]{babel}` line, it is not possible to pass options to the individual `.ldf` files directly. However, \LaTeX provides the command `\PassOptionsToPackage`, which declares options for a package before the package is loaded. So for example, `\PassOptionsToPackage{a=b,c=d}{foo.bar}` appends `a=b,c=d` to the macro `\opt@foo.bar`. `magyar.ldf` examines `\opt@magyar.ldf`, so for example passing options with `\PassOptionsToPackage{titles=\enskip}{magyar.ldf}` forces the space in section headings between the section number and the section title to be `\enskip`. (For compatibility reasons, `magyar.ldf` also processes the contents of `\magyarOptions` as options. This is useful to make `magyar.ldf` work with plain \TeX .)

\TeX macro wizards may enjoy studying the option parsing code in `magyar.ldf`. The entry point of the routine is `\processOptions`, whose argument is a comma-separated option list of `<key>=<value>` pairs. The routine calls `\processOption{<key>}{<value>}` for each pair found. This code is shown in figure 1.

2.13 Default Option Sets

Since there are 51 load options in the current `magyar.ldf`, the user should not be forced to know all of them. Reasonable defaults are provided (namely, `defaults=over-1.4`), so novice users can simply proceed. Intermediate users can select one of the five defaults, and possibly change a few options they don't like in their preferred default, and only expert users will change many options individually.

The number of bytes loaded were measured in a recent version of `magyar.ldf`, totalling 177 353 bytes. Out of that 177 kB, 32 417 bytes were used for initialization and providing the load option support framework, and declaring the options for the five defaults. After that, 138 872 bytes were used for implementing features selected by the load options. In the descriptions below, the number of feature bytes skipped is listed. (The larger the number, the less of `magyar.ldf` is processed at load time.)

The default sets are:

1. `=over-1.4` (10 720 bytes not loaded) This is the default among the defaults. Its main goal is to make all documents with the previous version of `magyar.ldf` (1.4) compile with the new version and to provide emergency bugfixes to incompatibility problems caused by the old version. It introduces a few essential typographical changes which have little impact on line

```

\def\processOptions#1{\processOptions@a#1,\hfuzz,}
\def\processOptions@a#1,{%
  \if,\noexpand#1,% (1)
    \expandafter\processOptions@a
  \else
    \csname fi\endcsname% (2)
    \processOptions@b#1,=%
  \fi}
@gobble\iftrue
\def\processOptions@b#1#2=#3,#4\fi% (3) (4)
  \ifx\relax#4\relax
    \ifx#1\hfuzz% Terminator
      \expandafter\expandafter\expandafter
        \gobble% (5)
    \else
      \if,\noexpand#1% OnlySpace
      \else% MissingArg; #2 ends by comma
        \if=\noexpand#1\missingKey#2%
          \else
            \missingVal#1#2\fi
        \fi
      \fi
    \else% Normal
      \processOption{#1#2}{#3}%
    \fi
  \processOptions@a}
\def\missingKey#1,{\errmessage{Key missing
  for value: #1}}
\def\missingVal#1,{\errmessage{Value (=) missing
  for option: #1}}
\def\processOption#1#2{\typeout{Got option
  key=(#1) val=(#2)}}

```

Figure 1: Option parsing code. Comments: (1) ignores extra commas, detects them by testing whether #1 is empty; (2) this is a `\fi` when expanded, but doesn't count as a `\fi` when being skipped over because its surrounding condition is false. The real `\fi` won't be expanded, because it is parsed as the parameter terminator of `\processOptions@b`. (3) needs #1#2 instead of just #1, so TeX will ignore space tokens in front of #1. As a side effect, when the option begins with =, the = will be put into #1, so `\missingKey` can be reported. (4) There are four different cases in which `\processOptions@b` can be invoked. The exact case is determined by how the macro parameter text separates parameters. The cases are: *Normal case*: #1#2 is the key, #3 is the value, #4 is =,; *MissingArg case* (=value) is missing, or (key) is missing, but =value) is present): #1#2 is (key), or =value), #3 and #4 are empty; *Terminator case*: #1#2 is \hfuzz, #3 and #4 are empty; *OnlySpace case*: #1 is , #2, #3 and #4 are empty. (5) `\gobble` removes the call of `\processOptions@a` at the end of the macro, so the iteration is finished.

and page breaks; it disables big, eye-catching changes. It makes most new commands available, but doesn't turn new features on.

2. `=compat-1.4` (82 890 bytes not loaded) Implements a strict compatibility mode with version 1.4 of `magyar.ldf`. Documents look about the same (exact match not guaranteed), even when the output is typographically incorrect. It does not define new commands such as `\told` or `\emitdate`.
3. `=safest` (124 679 bytes not loaded) Turns off all features, reverts to L^AT_EX and Babel defaults in every respect. It is useful only for debugging purposes: if a document doesn't compile, but it compiles with `defaults=safest`, individual options can be turned on one-by-one to see which is causing the compatibility problem.
4. `=prettiest` (1 221 bytes not loaded) Turns on all new features, and tries to follow Hungarian typography in the prettiest, most advanced way. It is possible that compatibility problems will arise with other packages, although the author is not currently aware of any.
5. `=hu-min` (1 317 bytes not loaded) Follows Hungarian typographical rules as closely as possible. Compliance is not complete, of course, because some aspects are not implemented; thus they are not covered by load options, and they cannot be controlled using defaults. If typographical rules allow choice (e.g. the first paragraph of a section may or may not be indented), the easiest and most compatible solution is chosen (e.g. accept the indentation defined by the document class).

2.14 Skipping Parts of the Input File

Since some parts of `magyar.ldf` can be disabled using load options, as we have seen, it is desirable to skip them completely. The easiest way of skipping part of TeX code is wrapping it into `\ifnum\MyFeature<1... \fi`. But this kind of skipping will consume hash memory for new control sequences skipped over, and it also requires that the skipped part is properly nested with respect to `\if...s`. `magyar.ldf` defines the following macro to do skipping without these flaws.

```

@gobble\iftrue
\def\skiplong#1{\fi
  \bgroup% so ^} would close it
  \catcode\string'^13
  \lccode\string'^=\string'^
  \lowercase{\let^~\fi}%
  \catcode\string'\14
  % comment, save hash memory

```

```

\catcode\string' $14
\iffalse}
\@gobble\fi
Now, code can be skipped with the construct
\ifnum\MyFeature<1 \skiplong\fi
...
\@gobble
{~}

```

`\lowercase` is needed in the implementation because `\let^~\fi` does not work, since the catcode of `~` is already assigned to be superscript when the definition of `\skiplong` is read.

2.15 Detecting Digits for Definite Articles

`\ref`, `\pageref` and `\cite` generate numbers, which are often prefixed by the definite article in Hungarian. The construct `'the \ref{foo}`; works fine in English, but the Hungarian definite article has two forms: *a* and *az*. *Az* must be used if the following words (as pronounced) starts with a vowel, and *a* must be used for consonants. So we need a macro that generates the definite article for numbers automatically. `magyar.ldf` contains the macro `\az`, which prefixes its argument by either `a~` or `az~`. This kind of macro is at present unique to `magyar.ldf`; other `.ldf` files apparently do not implement solutions for similar problems in other languages.

Unfortunately, `\az` is not expandable, because it redefines the meaning of several commands before processing its argument. `\az` works by half-expanding its argument (fully expanding, of course, `\ref`, `\pageref` and `\cite`), ignoring braces and most control sequences, non-digit and non-letter characters, changing `\romannumeral` to `\number` (so that *x* will become *az x*, but `\az{\romannumeral 10}` yields *a x*), looking for a number, a word or a single letter in the beginning (in fact, *az* has to be emitted if the starting digit is 5, and *a* has to be emitted if the starting digit is not 1 or 5). For words, single letters and positive numbers not beginning with 1, the proper definite article depends on the first letter only.

For numbers starting with 1, the definite article must be *az* if and only if the number of digits is $3k + 1$ for an integer *k*. For example, the Hungarian words for 1, 12, 123, 1000 are *egy*, *tizenkettő*, *százhuszonhárom*, *ezer*, respectively, and the definite forms are *az 1*, *a 12*, *a 123*, and *az 1000*, respectively. So we have to count the number of digits of a number.

It is not necessary to have 10 `\if` commands to test whether macro argument `#1` is a digit: it almost always works fine to use `\ifnum1<1\string#1`. The space at the end is important, because it will

terminate the second number of the `\ifnum` if `#1` is a digit. The condition $(1 < 1)$ is false if `#1` is a non-digit, and true $(1 < 10, 1 < 11 \text{ etc.})$ otherwise. `\string` cancels the special catcode `#1` might have. `#1` shouldn't be longer than a single token, because the test makes `TEX` process the extra tokens when `#1` contains a digit followed by extra tokens. If `#1` is `\if` or similar, the test isn't skippable. The test works even if `#1` is empty.

The macro `\Az` is also defined to insert a capitalized definite article at the beginning of a sentence. The macros `\aref`, `\Aref`, `\apageref`, `\Apageref`, `\acite` and `\Acite` are combinations of `\az` and referencing commands, so for example `\aref{foo}` is equivalent to `\az{\ref{foo}}`.

2.16 Counting Digits with Multiple Sentinels

A "sentinel" is something placed at the end of a list so that a conditional iteration over the list stops at the sentinel. For example, section 2.12 uses `\hfuzz` as a sentinel for the option processing of the macro `\processOptions@b`. A sentinel is usually a single token, but sometimes multiple sentinels have to be used in a row, when a macro processing them takes multiple parameters.

As mentioned in section 2.15, the Hungarian definite article (*a/az*) for a number depends on its pronunciation. The rule is: *az* has to be emitted for numbers starting with 5, and for numbers starting with 1 and having the number of digits following 1 divisible by 3. All other numbers are preceded by *a*. `magyar.ldf` thus contains a macro that counts number of digits following it:

```

\def\digitthree#1{\digitthree@#1//\hbox$}
\def\digitthree@#1#2#3{%
  \csname digitthree@%
    \ifnum9<1\string#1 \ifnum9<1\string#2
      \ifnum9<1\string#3 %
    \else b\fi\else b\fi\else z\fi\endcsname}
\def\digitthree@z#1\hbox${z}
\def\digitthree@b#1\hbox${}
\message{1:\digitthree{ } 100:\digitthree{23+}
  1000:\digitthree{456}}

```

In this example the macro `\digitthree` expands to *z* if its argument starts with digits of the multiple of 3. `\hbox$` is used as a sentinel to skip everything after the last digit has been found. The sentinel must not be present in the parameter itself. `\hbox$` makes no sense in `TEX`, so it is quite reasonable to assume that the parameter doesn't contain this. `magyar.ldf` uses `\hfuzz` and `\vfuzz` when only

a single token is allowed, because these two unexpandable tokens are quite rare. The three consecutive slashes in the example are three sentinels, so `\digitthree@` has always enough arguments.

However, the test doesn't work if the parameter of `\digitthree` contains braces. For example `\digitthree{1{2x}}` would look for the undefined control sequence `\digitthree@x`.

In the example the `\csname` trick was used to avoid `\expandafter` in the nested `\ifs`. See the definition of `\@@magyar@az@set` in `magyar.ldf` for using three multi-character sentinels in the same macro.

2.17 Definite Articles Before Roman Numerals

`\az` in `magyar.ldf` works differently for `\az{x}` and `\az{\romannumeral 10}` (see section 2.15), but how should it distinguish when `\romannumeral` has already been expanded by the time `\az` is called? Although there is no general solution to the problem, `magyar.ldf` addresses the case when `\az{\ref{my-part}}` is called, having the label `my-part` point to a `\part`, when `\def\thepart{\@Roman\c@part}` is active. (This is so with the standard `book.cls`.) `\ref` gets the part number from the `\newlabel{my-part}{x}{42}` command written to the `.aux` in the previous run of \LaTeX . `\label`, which has emitted this `\newlabel` has already expanded `\romannumeral` in the previous run, long before our `\ref` is called.

To make the definite article work in this special case, `magyar.ldf` redefines `\label` so it writes `\hunnewlabel` in addition to `\newlabel` to the `.aux` file. The arguments of `\hunnewlabel` are pre-expanded when `\let\romannumeral\number` is in effect. This solution also works when `\pageref` refers to a roman numeral page number.

Expanding the page number at the right time is rather tricky. The \LaTeX `\protected@write` says `\let\thepage\relax`, which prevents expansion in the following `\edef`, so `\thepage` is expanded only when the page is shipped out, and `\c@page` contains the right page number. What we want is to half-expand `\thepage`, so it gets expanded to `\@roman\c@page`, and `\@roman` is expanded to `\number (!)`, but the expansion of `\number\c@page` is postponed until the page is shipped out. This can be done by defining `\def\romannumeral{\noexpand\number}` before calling `\protected@write`. In practice, `magyar.ldf` itself expands the page number, so three `\noexpands` are needed in front of `\number`.

Redefining `\label` (so it emits `\hunnewlabel`) also raises a problem. Some packages loaded later might also override `\label`, for example `hyperref.sty` loads `nameref.sty` `\AtBeginDocument`, which over-

rides `\label`. `magyar.ldf` recognises the new definition when the Hungarian language is activated — which is done after the `\AtBeginDocument` hooks are run (see section 2.21). So `\hunnewlabel` works fine with `hyperref.sty`.

2.18 Removing All Braces

Removing all braces from a token list is required by the `\az` command (that inserts the *a/az* definite article). `\az` can find the first letter of its argument more easily if the argument doesn't contain braces.

The `\removebraces` macro defined in figure 2 removes all braces and spaces (recursively) from the tokens following it, until the first `\hfuzz`. The tokens may not contain a `\hfuzz` inside braces, but they may contain expandable material, even with unbalanced conditionals, because those are left unexpanded in `\removebraces@nobody` by `\noexpand`. The most important trick here is the construct `\ifcat{\noexpand#1`, which is true if `#1` starts with a brace, and yields `#1` with its first brace stripped. `\iffalse}\fi` is needed so that the macro definition is nested with respect to braces. The usage of `\@firstoftwo` is also worth mentioning: it is used to change the `\removebraces@nobody` token following the `\if` to `\removebraces`.

2.19 Changing \catcodes Safely

`\makeatletter` is equivalent to `\catcode 64 11` on ASCII systems; this changes the category code of characters having code 64 to 11 (letter). It is possible to specify the character `@` without knowing its character code: `\catcode'@12`. Wherever \TeX looks for a number (after `\catcode`, `\ifnum`, `\number`, etc.), it accepts a decimal number, an octal number prefixed by `'`, a hexadecimal number with digits 0–9A–F prefixed by `"`, an internal counter (such as `\linepenalty`), a `\count` register (such as `\count42` or `\@listdepth`) or a character prefixed by `'`. The character can be specified as a character token, or as a single-character control sequence. It is wise to specify `{, }, %` and space as `\{, \}, \%` and `_`, respectively, so the whole construct is properly nested with respect to braces, and since `%` and space tokens would be ignored.

However, many Babel language modules (`.ldf` files) make the character `'` active (i.e. `\catcode 13`), so the definition of `'` in `\catcode'@12` gets expanded. The expansion can be prevented by using `\noexpand`, but `\noexpand'` yields `'_13`, which is wrong, because `'_12` is needed, and moreover, will be expanded in the second run, because \TeX is looking for a number. Fortunately, `\string'` solves the problem, because `\string` changes the `\catcode` of

```

\@gobble\iftrue
\def\removebraces@stop#1#2\fi{#1}%
\def\removebraces#1{\ifx\hfuzz#1\removebraces@stop\fi
  \expandafter\removebraces\expandafter{\ifcat{\noexpand#1\hfuzz\iffalse}\fi
  \expandafter\removebraces\else\hfuzz}\removebraces@nob{#1}\fi}
\def\removebraces@nob#1#2{#2\ifx\hfuzz#1\hfuzz\expandafter\@firstoftwo% #2 is \fi
  \expandafter\removebraces\fi\removebraces@nobone#1}
\def\removebraces@nobone#1{\noexpand#1\removebraces}
\message{R:\removebraces {{foo}}{b}{{{a\fi}}}{r}}\hfuzz;}

```

Figure 2: `\removebraces`: A macro to remove all braces.

the following character token to 12 (other) or 10 (space); and, if a control sequence follows, `\string` converts it to a series of character tokens with `\catcode` other or space.

Thus, the ideal definition of `\makeatletter` is `\catcode\string'\@11␣`, which doesn't rely on the previous `\catcode` of `'` or of `@`. The space at the end of the definition is needed so `TEX` knows that the number 11 won't be followed by subsequent digits. Of course, the definition works only when the characters `catcodestring` have `\catcode` letter, `\` is an escape character (`\catcode` 0), and space is a space (`\catcode` 10). These are reasonable assumptions, because none of the standard `LATEX` packages change them.

The `LATEX` kernel's definition of `\makeatletter` is `\catcode'\@11\relax`, having `\relax` instead of space, which is equally good to mark the end of a number. This definition doesn't need `\string`, because at the time it is read, the `\catcode` of `'` is guaranteed to be 12 (other).

`magyar.ldf` saves the `\catcode` of `' ! & + - = | ; : ' " ? /` in the beginning, changes them to other, and restores them just before `\endinput`. This is needed in case other `.ldf`s have been loaded (e.g. `\usepackage[french,magyar]{babel}`) that have redefined `\catcodes`. For example, `french.ldf` activates `! ? ; :`.

It is also good not to change `\catcodes` until `\begin{document}` (not even `\AtBeginDocument`), because other packages not yet loaded may depend on the old, unchanged `\catcodes`. Babel, unfortunately, activates a character immediately when a shorthand is defined in an `.ldf` file, so this can raise strange compatibility issues—which can be partly resolved by loading most other packages before Babel. `magyar.ldf` solves this by not touching the `\catcode` of its own shorthands at the time of definition, but instead calls `\bbl@activate` in `\extrasmagyar`, and `\bbl@deactivate` in `\noextrasmagyar`. This is a local and temporary solu-

tion only. Future versions of Babel are expected to postpone character activation as far as `\@preamble` `cmds` (see also section 2.21).

2.20 Shorthands

A shorthand is an active character defined by an `.ldf` file with the `\declare@shorthand` command provided by Babel. In this sense, all active punctuation characters (see section 2.9) are shorthands.

The most important shorthand in `magyar.ldf` is `'13`. (Most `.ldf` files choose that character to be the main shorthand, but some, such as `germanb.ldf`, choose `"13`.) The use of Hungarian shorthands can be disabled by the `active=` load option, and the shorthand character can be changed from `'` with the `activeprefix=` load option. `magyar.ldf` also provides the `\shu` command, which is a longer form of `'13`, but without the possibly hazardous `\catcode` change.

Each shorthand is an active character, which raises compatibility problems (see section 2.19). `magyar.ldf` tries as hard as possible to avoid problems, but all efforts are in vain if another `.ldf` file is loaded which activates the same shorthand in the default (and unsafe) way.

For the user a shorthand is a control sequence without a backslash, so a shorthand is a command that can be typed and read quickly. `germanb.sty` provides `"a` to be equivalent to `\"a`, saving a keystroke for every accented German letter. `magyar.ldf` doesn't provide this saving, because the letters `o` and `u` have 3 accented forms, and introducing different letters for them would lead to confusion. Hungarian `LATEX` authors are encouraged to use the `latin2` encoding to type accented letters as a single character.

But the shorthand does an important job concerning (unaccented) long double consonants; for example, `'tty` is an abbreviation for `t\nobreak\discretionary{y-}{ }{ }ty\nobreak\hskip\z@skip`. (Section 2.6 explains why this is needed.) It should

be noted that shorthands are implemented as TeX macros, so ‘`_{}tj`’ and ‘`ttj`’ are equivalent.

The shorthand functionality of `magyar.ldf` for non-letters is inspired by `ukraineb.ldf`. ‘`=`’ and ‘`-`’ stand for a hyphen that separates words, so both words are automatically hyphenated by TeX (implemented as `\leavevmode\nobreak-\hskip\z@skip`); ‘`-`’ in math mode stands for a space character following a delimiter (`\mskip2.4mu plus3.6mu minus1.8mu`) that will magically be exactly as wide as if a space was inserted outside math mode, because the implicit `\mskip0.6mu` after the delimiter is already subtracted; ‘`--`’ emits `\,--\,`, to be used between author names in Hungarian bibliographies; ‘`|`’ emits a hyphen that is repeated at the beginning of the next line if the line is broken there (implementation: `\leavevmode\nobreak-\discretionary {}{-}{-}\nobreak\hskip\z@skip`), to be used with long words (e.g. *nátrium--klorid*) having important hyphens; ‘`_`’ inserts a discretionary hyphen with automatic hyphenation enabled at both sides; ‘`<`’ inserts a French opening guillemet even if the ligature `<<` is missing from the current font; ‘`>`’ inserts its paired closing; ‘`"`’ is equivalent to `\allowbreak` with hyphenation enabled on both sides (implementation: `\hskip\z@skip`); ‘`~`’ inserts a hyphen that doesn’t form ligatures when repeated (implementation: `\textormath{\leavevmode\hbox{-}}{-}`).

2.21 Inserting Code at `\@preamblecmds`

Babel calls `\selectlanguage` to set the default language `\AtBeginDocument`, which is (in general) too early. Suppose that the default language redefines `\catcodes` to be used with active characters. All packages that are loaded after the default language is activated will contain characters with unexpected and invalid catcodes. For example, if `hyperref.sty` is loaded after `magyar.ldf`, the `\AtBeginDocument` entries of `hyperref.sty` contain `\RequirePackage{name ref}`, which is executed *after* the entry `\selectlanguage{magyar}` of `babel.sty`, so `nameref.sty` will be loaded with wrong catcodes, and it will fail.

The solution is to postpone activation of the default language until after the `\AtBeginDocument` hooks. To accomplish this, `magyar.ldf` appends to `\@preamblecmds`, which is executed by the L^AT_EX kernel in `\document`, after `\AtBeginDocument`.

But what about the call to `\selectlanguage` inserted `\AtBeginDocument` by `babel.def`? Fortunately, it becomes a no-op, because `magyar.ldf` modifies `\selectlanguage` to do nothing if `\languagename` hasn’t changed — and this is exactly the case when activating the default language. On the other hand, `\@preamblecmds` has to force

the change even when `\languagename` is unchanged, so it calls `\select@language` (notice the at-sign). So `magyar.ldf` adds a call to `\select@language to \@preamblecmds`.

It also runs `\pagestyle{headings}` for the relevant document classes, so `\ps@headings` is executed once more, and the Hungarian version of the headings as defined by `magyar.ldf` will have a chance to be installed.

2.22 Displaying Theorem Titles

In English, theorem titles are displayed as “Theorem 1”, but Hungarian requires “1. tétel.”. To implement this, the `\@begintheorem` and `\@opargbegintheorem` macros are redefined each time the Hungarian language is activated. However, if `theorem.sty` or `ntheorem.sty` is loaded, the changes have to be embedded into a theorem style. The chosen name for the style is `magyar-plain`. It is activated by default when `magyar.ldf` is loaded, so theorem titles will come out right unless the user calls `\theoremstyle`. When `amsthm.sty` is loaded, `magyar.ldf` redefines the macros `\thmhead` and `\swappedhead` so both will emit the title properly.

2.23 Extra Symbols, `\paragraph` Titles, and Description Items

Hungarian typography requires a separator character other than a dot after the `\paragraph` title. Thus, a paragraph in English starting with “title text” should be something like “title ♦ text”. `magyar.ldf` provides several pre-defined title separation symbols, selected by the load option `postpara=`; similarly, `postsubpara=` controls `\subparagraphs` and `postdescription=` controls `\items` in the description environment.

`magyar.ldf` redefines `\paragraph` in a truly ugly fashion when `postpara=` is active, so that no extra horizontal space is inserted after the title, but the title ends at the separation symbol. The default definition of `\paragraph` is based on `\@startsection`, whose argument #5 is a negative skip, which means a positive horizontal skip after the title. This is changed to `-1sp` by `magyar.ldf` to avoid the skip, and an optional argument is always passed to the original `\paragraph` so the title will be typeset with the separator.

2.24 Indentation after Section Titles

Hungarian typography allows the first paragraph after a section title to be either indented or unindented, so `magyar.ldf` provides `afterindent=` as the load option to control this. L^AT_EX calculates the value of a boolean variable `\if@afterindent` from


```

\expandafter\addto\csname
  \expandafter\ifx\csname mathoptions@on
  \endcsname\relax check@mathfonts
  \else mathoptions@on\fi
\endcsname{\catcode'\,12 \mathcode'\,"8000
\begingroup\lccode'\, \lowercase
{\endgroup\def~}{\HuComma}}

```

Figure 3: `\HuComma`: Smart commas in math.

the signedness of a parameter of `\@startsection`, and later uses that boolean to insert or omit the indentation. The value is forced to true by `magyar.ldf` by the simple definition `\let\@afterindentfalse \@afterindenttrue`.

2.25 The Decimal Comma

The dot character is defined as *ordinary* in mathematical text by default, so decimal real numbers can be typed simply as `--12.34`. Hungarian denotes the decimal point by a comma instead of a dot, but typing `-$-12,34$` yields ‘-12,34’ with too much space after the comma, because the comma is defined as *punctuation* rather *ordinary* in math text. `-$-12{,}34$` yields ‘-12,34’, which is correct, but `magyar.ldf` provides two mechanisms to save the two keystrokes of the curly braces around the comma.

First, the `\HuComma` macro below inserts an ordinary comma if it is followed by a digit, and an operator comma otherwise:

```

\edef\hucomma@lowa#1#2 #3#4 #5#6\hfuzz{%
  \noexpand\ifnum9<1#5 \noexpand\if#1t%
  \noexpand\if#3c% (1)
  \noexpand\mathord\noexpand\fi\noexpand\fi%
  \noexpand\fi\mathchar%
  \ifnum\mathcode'\,"8000 "613B \else\the%
  \mathcode'\, \space\fi}%
\def\hucomma@lowb{\expandafter\hucomma@lowa
  \meaning\reserved@a / \hfuzz}%
\DeclareRobustCommand\HuComma
  {\futurelet\reserved@a\hucomma@lowb}

```

The solution Donald Arseneau proposed to the comma problem inspired these macros. In line (1) `\hucomma@lowa` tests whether the `\meaning` of the following character is ‘the character <digit>’. A `\meaning` is always at least three words, but it may be more (e.g. ‘math shift character \$’). Only *the character* starts with letters *t* and *c*. An `\edef` is needed above so the `\mathchar` emitted doesn’t depend on `\mathcode` changes after the definition of `\HuComma`. Then the comma character can be redefined as `\HuComma`, as given in figure 3.

With these definitions, the formula ‘ $F_i(x, y) = y^i + 1,3x$, $x, y \in A$, $i = 1, 2, 3, \dots$ ’ can be typed simply as `$F_{i}(x,y)=y^i+1,3x, \ x,y \in A`,

`\ i=1, \ 2, \ 3, \ldots$`, if `_` is breakable (such as in `nath.sty`).

When `nath.sty` is loaded, the definitions are appended to `\mathoptions@on`, and if `nath.sty` is missing, to `\check@mathfonts`. The appropriate macro is run just before `\everymath` by L^AT_EX. Redefining the `\catcode` and `\mathcode` this way ensures that the proper comma is used inside math mode—unless the whole math formula is a macro argument with already assigned `\catcodes`. Also, it is not a good use of `\begingroup`, `\lccode`, `\lowercase` and `\endgroup` to modify the active meaning of a character without actually activating it. Calling `\catcode'\,13` before `\def` wouldn’t help here anyway if the whole construct is embedded into a macro definition, because `\catcode` wouldn’t be able to change an already assigned `catcode`.

`frenchb.ldf` provides `\DecimalMathComma` and `\StandardMathComma` to change the `\mathcode` of the comma. However, the smart comma based on `\HuComma` acts correctly without the user needing to be aware of curly braces or redefinitions.

The solution above can be activated with the loading option `mathhucomma=fix`. An alternative approach doesn’t alter `\mathcodes`, but introduces a special math mode in which the dot appears as a comma only when the Hungarian language is active. Thus the printout of `\MathReal{-12.34}` depends on the current Babel language. The definition of `\MathReal` in `magyar.ldf` is similar to:

```

\def\mathreal@lowa#1{\ensuremath{%
  \mathreal@lowb#1\@gobble.}}
\def\mathreal@lowb#1.{%
  #1\@secondoftwo\@gobble(1)
  {\mathchar"013B \mathreal@lowb}}% comma
\DeclareRobustCommand\MathReal{\ensuremath}%
{\catcode'\ 11\relax\addto\extrasmagyar{%
  \babel@save\MathReal %
  \let\MathReal \mathreal@lowa}}

```

The argument of `\MathReal` must contain the dot to be changed literally, outside braces. There is a little macro wizardry in the implementation that stops calling `\mathreal@lowb` infinitely. The call `\mathreal@lowa` terminates its argument by a sentinel `\@gobble.`, so `#1` of `\mathreal@lowb` will end by `\@gobble`, which will gobble `\@secondoftwo`, so the `\@gobble` in line (1) will take effect, which stops the recursion.

`\MathReal` is going to be extended in the future so it will handle physical units following the number properly, and it will also insert thin spaces after each three digits. This feature has already been implemented in `frenchb.ldf`.

2.26 Parsing Dates

There are many correct ways to write dates in Hungarian, and `magyar.ldf` provides an `\emitdate` command that can generate any of these formats. Doing the reverse is a little more interesting.

Let's suppose we have a Gregorian date consisting of a year (4 or 2 digits), a month (a number or a name) and a day-of-month in some standard format. We want a command `\parsedate` to detect the format, split the date into fields, and call `\fixdate`:

```
\def\fixdate#1#2#3{%
  \@tempcnta#1 \ifnum#1<50
  \advance\@tempcnta2000 \fi
  \ifnum\@tempcnta<100
  \advance\@tempcnta1900 \fi
  \typeout{found year=(\@tempcnta)
  month=(#2) day=(#3)}%
}

\def\stripdot#1{\expandafter%
  \stripdot@lowb\stripdot@lowa
  #1\relax.\relax}%
\def\stripdot@lowa#1.\relax{#1\relax}%
\def\stripdot@lowb#1\relax#2\relax{#1}%
```

The definition of `\parsedate` is shown in figure 4. `\parsedate` first does some generic cleanup, and puts the resulting date into `\re@b`. `\endgroup` cancels the redefinition of `\today` etc., but `\re@b` is expanded first, which defines itself, so the value of `\re@b` will be retained after `\endgroup`. After that, the trailing dot is stripped, and then various `\parsedate@...` commands are run. If a command recognises the date format, it puts a call to `\fixdate` into `\re@a`, which will be called at the end of `\parsedate`. Strange strings like `!//:!\hfuzz` are sentinels.

The idiom `\expandafter\endgroup\re@b` is an important trick for expanding a macro before the current group completes (and changes are undone). It usually contains definitions of other control sequences whose meanings are about to be retained after the end of the group. An alternative would be to inject such a definition using `\aftergroup`, but that only accepts a single token, so it would be very painful to make a macro definition with spaces and braces survive this way.

The individual `\parsedate@...` commands are given in figure 5. This implementation of date parsing isn't error-proof. If something weird is passed to `\parsedate`, it may produce surprising TeX er-

```
\def\parsedate#1{%
  \begingroup
  \def\today{\the\year-\the\month-\the\day}%ISO
  \let\protect\string
  % remove accents from Hungarian month names:
  \let'\@firstofone
  \let~\space%change '2003.~okt' to '2003. okt'
  \edef\re@b{\def\noexpand\re@b{#1}}%
  \expandafter\endgroup\re@b
  \edef\re@b{\expandafter\stripdot\expandafter
    {\re@b}}%
  \let\re@a@empty \expandafter\parsedate@a\re@b
  !--!\hfuzz
  \ifx\re@a@empty \expandafter\parsedate@f\re@b
  !//:!\hfuzz \fi
  \ifx\re@a@empty \expandafter\parsedate@b\re@b
  !//!\hfuzz \fi
  \ifx\re@a@empty \expandafter\parsedate@c\re@b
  !..!\hfuzz \fi
  \ifx\re@a@empty \expandafter\parsedate@d\re@b
  !.xyz !\hfuzz \fi
  \ifx\re@a@empty \expandafter\parsedate@e\re@b
  !xyz , !\hfuzz \fi
  \ifx\re@a@empty \errmessage{Unrecognised date:
    \re@b}%
  \else \re@a% call \fixdate
  \fi}
```

Figure 4: `\parsedate`: Parse date formats.

rors. However, `\parsedate` can distinguish between different formats of correct input.

2.27 Setting Up French Spacing

Hungarian typography requires `\frenchspacing` to be turned on, but most L^AT_EX users fail to follow this requirement. Babel provides the command `\bbl@frenchspacing`, which turns French spacing on if it was off. The `frenchspacing=` load option of `magyar.ldf` controls how Hungarian text should behave. For the sake of symmetry, `magyar.ldf` provides `\@@magyar@antifrenchspacing`, which — contrary to the typographical requirement — turns french spacing off:

```
\def\@@magyar@antifrenchspacing{%
  \ifnum\the\sfcode'\.=\@m
  \nonfrenchspacing
  \let\@@magyar@nonfrenchspacing%
  \frenchspacing
  \else \let\@@magyar@nonfrenchspacing\relax
  \fi}
\let\@@magyar@nonantifrenchspacing%
  \frenchspacing
\addto\extrasmagyar{\@@magyar%
  @antifrenchspacing}
\addto\noextrasmagyar{\@@magyar%
```

```

\def\parsedate@a#1-#2-#3!#4\hfuzz{% ISO date: YYYY-MM-DD
  \ifx\hfuzz#4\hfuzz\else \ifnum1<1\string#1\relax
  \ifnum1<1\string#2\relax \ifnum1<1\string#3\relax
  \def\re@a{\fixdate{#1}{#2}{#3}}%
  \fi\fi\fi\fi}
\def\parsedate@b#1/#2/#3!#4\hfuzz{% LaTeX date: YYYY/MM/DD
  \ifx\hfuzz#4\hfuzz\else \ifnum1<1\string#1\relax
  \ifnum1<1\string#2\relax \ifnum1<1\string#3\relax
  \def\re@a{\fixdate{#1}{#2}{#3}}%
  \fi\fi\fi\fi}
\def\parsedate@c#1.#2.#3!#4\hfuzz{% English date: YYYY.DD.MM
  \ifx\hfuzz#4\hfuzz\else \ifnum1<1\string#1\relax
  \ifnum1<1\string#2\relax \ifnum1<1\string#3\relax
  \def\re@a{\fixdate{#1}{#3}{#2}}%
  \fi\fi\fi\fi}
% vvv English and Hungarian month names
\def\mon@jan{1} \def\mon@feb{2} \def\mon@mar{3} \def\mon@apr{4}
\def\mon@maj{5} \def\mon@may{5} \def\mon@jun{6} \def\mon@jul{7}
\def\mon@aug{8} \def\mon@sze{9} \def\mon@sep{9} \def\mon@okt{10}
\def\mon@oct{10} \def\mon@nov{11} \def\mon@dec{12}
\def\parsedate@d#1.#2#3#4#5 #6!#7\hfuzz{% {2003. oktober 25}
  \ifx\hfuzz#7\hfuzz\else
  % now: {#1}=={2003}, {#2#3#4#5}=={oktober}, {#6}=={25}
  \ifnum1<1\string#1\relax \ifnum1<1\string#6\relax
  \lowercase{%
    \expandafter\ifx\csname mon@#2#3#4\endcsname\relax\else
    \edef\re@a{\noexpand\fixdate{\number#1}{%
      \csname mon@#2#3#4\endcsname}{\number#6}}\fi}%
  \fi\fi\fi}
\def\parsedate@e#1#2#3#4 #5, #6!#7\hfuzz{% {October 25, 2003}
  \ifx\hfuzz#7\hfuzz\else
  \ifnum1<1\string#5\relax \ifnum1<1\string#6\relax
  \lowercase{%
    \expandafter\ifx\csname mon@#1#2#3\endcsname\relax\else
    \edef\re@a{\noexpand\fixdate{\number#6}{%
      \csname mon@#1#2#3\endcsname}{\number#5}}\fi}%
  \fi\fi\fi}
\def\parsedate@f#1/#2/#3:#4!#5\hfuzz{% LaTeX default \today
  % YYYY/MM/DD:XX:YY
  \ifx\hfuzz#5\hfuzz\else \ifnum1<1\string#1\relax
  \ifnum1<1\string#2\relax \ifnum1<1\string#3\relax
  \def\re@a{\fixdate{#1}{#2}{#3}}%
  \fi\fi\fi\fi}

```

Figure 5: Individual `\parsedate...` commands.

`@nonantifrenchspacing}`

2.28 varioref.sty Fixes

The `magyar` load option of `varioref.sty` (2001/09/04 v1.3c) is buggy, because it uses the never-defined `\aza` command for adding definite articles, and it also calls `\AtBeginDocument` too late, producing a \TeX error each time the Hungarian language is activated. `magyar.ldf` contains the correct definitions for the language-specific text reference macros, such as `\reftextlabelrange`, and also contains ugly fix-up code to remove the wrong macros inserted by `varioref.sty`. A patch has been sent recently to the author of `varioref.sty`.

Some of these text reference macros use the `\az` and the `\told` commands defined by `magyar.ldf`.

2.29 Removing Full Stops After Section Titles in AMS Classes

AMS document classes always append a full stop after section titles, which is strictly forbidden in Hungarian typography. The solution is to remove the tokens `\@addpunct.` from the definition of `\@sect` (and also from `\NR@sect` in case `nameref.sty` has also been loaded). But this simple idea is quite complicated to program, and the result is ugly, as seen in figure 6. This detects AMS classes by the presence of `\global\@nobraektrue\@xsect` in the definition of `\@sect`, and adds code just before `\@xsect`. The code added prepends `\let\@addpunct\@gobble` to the definition of `\@svsechd`. `\@svsechd` is later called by `\@xsect`, which calls `\@addpunct`, but by that time `\@addpunct` is a no-op. The application of this fix is controlled by the `amspostsectiondot=` load option.

2.30 Reduced Math Skips

Investigations in [1] have shown that the following settings produce the desired space for Hungarian math mode:

```
\thickmuskip 4mu plus 2mu minus4mu
% LaTeX: 5mu plus5mu
\medmuskip 2mu plus1.5mu minus2mu
% LaTeX: 4mu plus2mu minus4mu
\thinmuskip 3mu % LaTeX: ditto
```

Notice that `\medmuskip < \thinmuskip`. These settings can be selected in `magyar.ldf` with the load option `mathmuskips=`. The difference between the original and the reduced spacing:

$$a + b - c/d * y \circ x = z \quad a + b - c/d * y \circ x = z$$

2.31 Breaking a Long Inline Math Formula

Hungarian typography requires that a binary relation or operator (e.g. in $1 + 2 = 3 + 4$) must be repeated in the next line if an inline math formula is broken there. This can be accomplished for the equation sign by substituting `=\nobraek\discretionary-{}{\hbox{\(=\)}}{}` for `=` in math formulas. The long inline formula delimiters `\(` and `\)` are used because the catcode of the `$` would be wrong if `nath.sty` was loaded after `magyar.ldf`. `\nobraek` is necessary, so \TeX itself won't break the line after the `=`.

The `mathbrk=` load option of `magyar.ldf` controls whether the operators and relations should be redefined. If so, the operators `+`, `-`, `*` (as well as 37 operators available as control sequences) and the relations `<`, `>`, `=`, `:` (as well as 43 relations available as control sequences) are modified so they get repeated at the beginning of the line. The `\cdot` and the `\slash` operators are also modified, because Hungarian typography disallows breaking the line around them.

2.32 Restarting Footnote Numbering on Each Page

Although `\usepackage[perpage]{footmisc}` and `footnag.sty` provide these features, `magyar.ldf` allows normal arabic footnote and page-restarting asterisk-footnotes to be intermixed. It is common in Hungarian article collections to have the notes of the author numbered in arabic (by `\footnote`), and the footnotes of the editor added with asterisks (by `\editorfootnote`). The first four editorial footnotes on a page are marked with `*`, `**`, `***`, and `†`. `magyar.ldf` also inserts proper additional space between the footnote mark and the footnote text, and the footnote facility is fully customizable with the `\footnotestyle` command.

The basic idea behind the implementation of pagewise numbering is creating a `\label` for each footnote, and whenever the `\pageref` for that label shows a different page, resetting the counter to zero. This clobbering can be automated by abusing the `\cl@footnote` hook. Each time `\stepcounter` advances a counter, the corresponding `\cl@...` hook is called, which usually resets other counters (for example, advancing the chapter counter resets the section counter). But arbitrary code can be executed after the automatic `\stepcounter{footnote}` by appending that code to the macro `\cl@footnote`.

The famous problem of creating a macro that will expand to n asterisks is proposed in appendix D of *The \TeX book* [3]. David Kastrup has provided

```

\expandafter\amssect@fixa\@sect [] [] [] [] [] [] [] %
  \global\@nobreaktrue\@xsect\hfuzz\@sect
\expandafter\amssect@fixa\@sect [] [] [] [] [] [] [] %
  \global\@nobreaktrue\@xsect\hfuzz\NR@sect% with nameref.sty
\long\def\amssect@fixa#1\global\@nobreaktrue\@xsect#2\hfuzz#3{%
  \ifx\hfuzz#2\hfuzz\else\amssect@fixb#3\fi}% fix if found
\def\amssect@fixb#1{% #1 is \@sect or \NR@sect
  \expandafter\let\csname amssect@saved\string#1\endcsname#1%
  \edef#1{\noexpand\expandafter\noexpand\amssect@low\expandafter
    \noexpand\csname amssect@saved\string#1\endcsname}%
  \let\@svsechd\@empty% prevent Undefined \cs
\long\def\amssect@low##1\global\@nobreaktrue{##1%
  \expandafter\def\expandafter\@svsechd\expandafter{%
    \expandafter\let\expandafter\@addpunct\expandafter\@gobble
    \@svsechd}%
  \global\@nobreaktrue}}

```

Figure 6: Removing full stops after AMS section titles.

a brilliant solution to the problem in [2], namely `\expandafter\mtostar\romannumeral\number{n}000A`, where `\mtostar` transforms `ms` to asterisks: `\def\mtostar#1{\if#1m*\expandafter\mtostar\fi}`. This solution is used in `magyar.ldf`.

`magyar.ldf` also provides the following command to insert footnotes into section titles such that neither the table of contents nor the page headings are affected:

```

\def\headingfootnote{%
  \ifx\protect\@typeset@protect%
    \expandafter\footnote
  \else\expandafter\@gobble\fi

```

2.33 Class-specific Modifications

`magyar.ldf` does some modifications based on the current document class (using the `\@ifclassloaded` \LaTeX command). Only the standard classes `article.cls`, `report.cls`, `book.cls` and `letter.cls` are supported at present. The visual appearance of the `\part` and `\chapter` output is changed, and the page headers are also modified. For `book.cls`, part numbering is spelled out, so “Part 1” becomes “Első rész” (“Part One”) if the load option `partnumber=` is set to `Huordinal`.

The command `\ps@headings` has to be executed again to install its changed heading macros. This is called from `\@preamblecmds`, after the default language has been activated (see section 2.21).

The typographically correct customization of `letter.cls` is under development.

2.34 Spelling Out Numerals and Ordinals

The `\hunumeral` and `\huordinal` macros defined in `magyar.ldf` can spell out integers between

–9999 and 9999. `\@Hunumeral` and `\@Huordinal` are the capitalized versions of these macros. For example, `\@huordinal{2004}` produces *kétezzer-nygyedik* (“two thousand and fourth”) and `\@Hunumeral{2004}` produces *Kétezzer-nygy* (“Two thousand and four”). All of these macros are fully expandable, so they can be used for `\part` numbering: `\def\thepart{\@Huordinal\c@part}`, or more simply: `\def\thepart{\Huordinal{part}}`.

The most important implementation issue is the method to retrieve the last digit of a number in an expandable construct. If the number is between 0 and 9999, the following macro solves the problem:

```

\def\LastDigitOf#1{\expandafter%
  \lastdigit@a\number#1;}%
\def\lastdigit@a#1;{% #1 in 0..9999
  \ifnum#1<10 #1\else\ifnum#1<100
    \lastdigit@b00#1%
  \else\ifnum#1<1000 \lastdigit@b0#1%
    \else\lastdigit@b#1\fi
  \def\lastdigit@b#1#2#3#4{#4}

```

2.35 Suffix Generation

As mentioned earlier, the Hungarian language has suffixes to represent relations in space and time, instead of prepositions. For example, an English math text might contain “It follows from (1)”, in which “from (1)” can be typed as `\ref{eq1}`. The \LaTeX referencing scheme guarantees that the text above will come out right, even if the order of equations is changed in the document.

But in Hungarian, the suffix standing in place of “from” has two forms: *-ból/-ből*, depending on the vowel harmony of the pronunciation of `\ref{eq1}`. So there is a need for automatic suffix generation.

`magyar.ldf` provides the command `\told`, with which the Hungarian version of “from (1)” can be typed as `\told(\ref{eq1})+bol{}`, which will generate “(1)-ból”, “(2)-ből”, but “(3)-ból”.

`\told` can handle 20 different suffixes, and 4·20 suffix combinations (such as `\told3+adik+ra{}`, meaning “to the third”). Only the last number is considered in references containing multiple numbers. Roman numerals are recognised properly in references with the help of `\hunnewlabel` (see section 2.17—this is implemented similar to `\az`). Suffix generation is supported only for integers and Hungarian document structure names (see section 3.2), because writing a generic suffix generator without a database is quite a difficult task, and definitely won’t give Hungarian L^AT_EX users much more comfort beyond the current `\told` implementation.

Although most Hungarian suffixes have 1, 2 or 3 forms,³ numbers can be classified into 23 paradigm classes, so that the paradigm class uniquely determines the correct form of all known suffixes. The reason that there are so many classes is because the letter *v* of the *-val/-vel* suffix must be also changed to the last letter of the number if that letter is a consonant. Essentially each final digit has a class, and there are classes for the powers of 10, and some of the numbers 20, 30, . . . 90 also have their own classes. To sum up, the suffix of a number depends on the last nonzero digit, and the number of trailing zeroes.

The implementation of `\told` is surprisingly long and ugly, full of recursive macros that parse the input, and it doesn’t contain any bright ideas that are not also found elsewhere in `magyar.ldf`. The curious T_EX hacker should study `\az` instead, because it is shorter and its trick density is much higher.

2.36 Warning Messages

`magyar.ldf` has the unique feature that it displays warning messages (called ‘suggestions’) at load time to notify the user that they are using `magyar.ldf` in a possibly incorrect way. If they are not disabled by the `suggestions=` load option, the following suggestions are displayed to standard output during the `\AtBeginDocument` hook:

- the user forgot to load `\usepackage{t1enc}`—so words with accented letters won’t hyphenate automatically;
- the user forgot to load `\usepackage[latin2]{inputenc}`, or the input encoding chosen is not `latin2`, `cp1250` or `utf8`—so there is a

³ Of course, suffixes with only one form are not supported by `\told`.

good chance that accented characters will disappear or come out wrong;

- the Hungarian hyphenation patterns requested were not found—`magyar.ldf` tries to use the other two possible Hungarian patterns, if they are available;
- `\def\magyarOptions` or `\PassOptionsToPackage{...}{magyar.ldf}` was specified too late—late options can be detected, but they have no effect, since options do their work while `magyar.ldf` is being loaded;
- the buggy `varioref.sty` has been loaded as `\usepackage[magyar]{varioref}`—this will happen until the patch is integrated to `varioref.sty`; the current version is so buggy that it displays an untraceable L^AT_EX error each time the Hungarian language is activated (see section 2.28).

3 Miscellaneous Tricks

First we show some common expansion tool macros defined by L^AT_EX:

```
\def\@empty{}
\long\def\@gobble#1{}
\long\def\@gobbletwo#1#2{}
\long\def\@firstofone#1{#1}
\long\def\@firstoftwo#1#2{#1}
\long\def\@secondoftwo#1#2{#2}
```

`\@firstofone` differs from `\@empty`, because it may not be followed by `}`, it ignores spaces in front of its argument, and it removes at most one pair of braces around its argument. All of these properties are consequences of the macro expansion rules described in chapter 20 of *The T_EXbook* [3].

This remainder of this section describes T_EX macro and typesetting tricks not tightly related to the Hungarian language.

3.1 The Factorial Sign in Math Mode

`nath.sty` contains a smart definition of the factorial operator, so $(a + b)!/a!b! + c! \cdot d!$, with proper spacing can be typed as `$(a+b)!{/a!b!+c!}\cdot d!$`. The only place where braces are needed is before the slash. `magyar.ldf` adapts the definition:

```
\def\factorial{\mathchar"5021\mathopen{}}%
\mathinner{}}
\expandafter\addto\csname \expandafter\ifx
\csname mathoptions@on\endcsname\relax
% detect nath.sty
check@mathfonts\else mathoptions@on\fi
\endcsname{\catcode'\!12
\mathcode'\!8000
\begingroup\lccode'\~!\lowercase{%
\endgroup\def~}{\factorial}}
```

```

\def\NumbersToOne#1{\nonumbers@a#1\hfuzz}
\def\nonumbers@skipa#1\nonumbers@s#{#1\nonumbers@a}
\def\nonumbers@a#1{% change first digit
  \ifx#1\hfuzz \expandafter\@gobble
  \else\if1<1\string#1\else\if\noexpand#1mi%
  \else\if\noexpand#1di\else\if\noexpand#1ci%
  \else\if\noexpand#1li\else\if\noexpand#1xi%
  \else\if\noexpand#1vi\else\if\noexpand#1ii%
  \else\if\noexpand#1MI\else\if\noexpand#1DI%
  \else\if\noexpand#1CT\else\if\noexpand#1LI%
  \else\if\noexpand#1XI\else\if\noexpand#1VI%
  \else\if\noexpand#1II%
  \else\noexpand#1\nonumbers@skipa
  \fi\fi\fi\fi\fi\fi\fi\fi\fi\fi
  \fi\fi\fi\fi\fi \nonumbers@s}
\def\nonumbers@s#1{% gobble next digits
  \ifx#1\hfuzz \expandafter\@gobble
  \else\if1<1\string#1\else\if\noexpand#1m%
  \else\if\noexpand#1d\else\if\noexpand#1c%
  \else\if\noexpand#1l\else\if\noexpand#1x%
  \else\if\noexpand#1v\else\if\noexpand#1i%
  \else\if\noexpand#1M\else\if\noexpand#1D%
  \else\if\noexpand#1C\else\if\noexpand#1L%
  \else\if\noexpand#1X\else\if\noexpand#1V%
  \else\if\noexpand#1I%
  \else\noexpand#1\nonumbers@skipa
  \fi\fi\fi\fi\fi\fi\fi\fi\fi\fi
  \fi\fi\fi\fi\fi \nonumbers@s}

```

Figure 7: `\nonumbers`: Change all roman and arabic numbers to 1.

3.2 Including the Structure Name in References

Text like “in subsection 5.6” is usually typed as `in subsection~\ref{that}`. But it would be nice if \LaTeX were able to guess that `\ref{that}` actually points to a subsection. The structure depth can be deduced by counting the dots in the reference: a subsection has one dot (in an article), and a subsection has two dots.

`magyar.ldf` provides the `\refstruc` command which has a smarter detection scheme: it changes all roman and arabic numbers to one (1, i or I) in the reference, and compares the result with the tokens generated by `\thechapter`, `\thesection` etc., with `\c@chapter` etc. set to 1 temporarily. This should work in most cases, although it cannot refer to equations, tables or figures. In Hungarian text, the Hungarian structure names are emitted, and other texts the original, English control sequence names are printed. `\refstruc` includes the definite article and suffixes support; for example,

```
\Az{\refstruc{that+tol}}
```

may emit *az 1. fejezettől* (“from chapter 1”).

The full implementation is quite long, and is not included here, but the macro that changes all roman and arabic numbers to one is presented in figure 7.

```

\ifnum 0<%
  \expandafter\ifx\csname setTrue\endcsname
  \relax\else\fi
  \expandafter\ifx\csname allowttyhyphens\endcsname
  \relax\else\fi
\space
\def\amsfix#1#2#3#4#5#6#7\vfuzz{%
  \def\reserved@a{#6}%
  \ifx\tocline#2\ifx\reserved@a\@empty%
    \def#1{\@tocline{#3}{#4}{#5}{\parindent}{}}%
  \fi\fi
  \expandafter\amsfix\expandafter\l@table\l@table
  ,,,,,,\vfuzz
  \expandafter\amsfix\expandafter\l@figure\l@figure
  ,,,,,,\vfuzz
\fi

```

Figure 8: Fixing overfull `\hboxes` in AMS classes.

3.3 Enabling Long Page Numbers

If the width of a page number in the table of contents is greater than `\@pnumwidth`, \LaTeX emits an “Overfull `\hbox`” warning. This can be eliminated by changing `\@dottedtocline` in the \LaTeX kernel. The line `\hb@xt@{\@pnumwidth}{\hfil \normalfont \normalcolor #5}` should be changed to:

```

\setbox\@tempboxa\hbox{\normalfont
R \normalcolor #5}%
\ifdim\wd\@tempboxa<\@pnumwidth\setbox%
\@tempboxa\hb@xt@{\@pnumwidth}{\hfil\unhbox
\@tempboxa}\fi \box\@tempboxa

```

Although this change isn’t related to the Hungarian language, `magyar.ldf` will do it given the `dottedtocline=` load option.

3.4 Removing AMS Warnings from \listoftables

Some AMS document classes (such as `amsart.cls`) produce an “Overfull `\hbox`” warning for each line in the `\listoftables` and `\listoffigure`. This can be fixed by changing this line in the `\l@table` and `\l@figure` macros in the AMS classes:

```

\@tocline{0}{3pt plus2pt}{0pt}{-}{-}
to:
\@tocline{0}{3pt plus2pt}{0pt}{-}{\parindent}{-}.

```

The code shown in figure 8 makes this change.

The control sequences `\allowttyhyphens` and `\setTrue` are defined by each of the AMS document classes, so their presence indicates that one of those classes are loaded. The logical or operation using `\ifxs` nested to the `\ifnum` test is also worth noting.

3.5 Discarding to End of File

The naïve solution `\ifskiprest\endinput\fi` results in the \TeX error message “`\end` occurred when `\iftrue` in line *n* was incomplete” if there is a line

break at the ¶ sign. Without the line break, the naïve solution works perfectly, because `\endinput` stops reading the current file *after* the current line, so the `\fi` also gets evaluated.

It is possible to do something before `\endinput`:

```
\expandafter\ifx\csname
  ver@foo.sty\endcsname\relax
  \endinput \errmessage{I am incompatible
    with foo.sty}\fi
```

All of the above must be put after `\endinput` without a line break. The `\csname fi\endcsname` construct closes the `\ifx` when the condition is true, but is invisible when the condition is false, and \TeX is skipping tokens.

The \LaTeX kernel macro `\ifpackageloaded` implements the conditional end by a different trick. The following two constructs are equivalent:

```
\ifpackageloaded{foo}{\endinput\errmessage{I
  am incompatible with foo.sty}}{}
```

and

```
\expandafter\ifx\csname ver@foo.sty\endcsname
  \relax\expandafter@gobble
\else \expandafter@firstoftwo \fi
{\endinput
\errmessage{I am incompatible with foo.sty}}
```

In the trick above, `\errmessage` isn't on the same line as `\endinput`. This isn't a problem, because by the time `\endinput` is evaluated by \TeX 's stomach, `\errmessage` has been read from the file, and it is already on the input stack. `\endinput` doesn't discard the input stack, it just prevents more lines from being read from the current file.

3.6 Typesetting Text Verbatim

`\catcodes` are assigned when \TeX 's eyes read the next character from the current line, so a `\catcode` command affects all subsequent characters of the current line, as well as the following lines. But once a category code has been assigned, it won't be affected by subsequent `\catcode` commands. For example, `@firstofone{\catcode'A 14 AAA}` makes `A` into a comment start character, but it emits three `As`, because the category code of the three `As` is already set by the time `\catcode` is executed.

The reason why the `\verb` command of the \LaTeX kernel cannot be part of a macro argument is the same: `\verb` changes the `\catcode` of most characters to *other*, but these changes have no effect inside a macro argument, because the argument has been read from the input file by the time `\catcode` can take effect.

Instead of altering `\catcodes`, a verbatim macro can be based on the `\meaning` primitive, so that it can be passed as an argument. However, \TeX 's

eyes will have destroyed some information such as comments and the exact number of successive spaces by the time `\meaning` is expanded. For example, these definitions are from `binhex.dtx`:

```
\def\verbatimize#1{\begingroup
  \toks0{#1}\edef\next{\the\toks0}%
  \dimen0\the\fontdimen2\font
    \fontdimen2\font=0pt
  \expandafter\stripit \meaning\next
  \fontdimen2\font=\dimen0 \endgroup}
\def\stripit#1>{}
```

3.7 Stopping the Iteration

Let's suppose we need a macro that capitalizes all *as* and *bs* until the first “.”:

```
\def\ucab#1{%
  \if\noexpand#1.\expandafter@gobble
  \else\if\noexpand#1aA%
  \else\if\noexpand#1bB%
  \else\noexpand#1%
  \fi\fi\fi\ucab}
\message{\ucab abc.abc} % -> ABcab
```

`\noexpand` prevents expansion of `#1` in case it is an expandable control sequence such as `\the` or a macro. If `\if` is changed to `\ifx`, then not only the character codes, but also the category codes would be compared.

The trick that stops the iteration here is that `\expandafter` expands the first `\else` that will remove everything up to the last `\ucab`. Then comes `@gobble`, which removes `\ucab`, and the iteration is stopped.

The construct doesn't work when `#1` has braces around it, or it is `\if...`, `\else` or `\fi`. Also, spaces will be ignored because of macro expansion.

But what if we'd like to capitalize only the *first a* or *b*? Then we would need `\expandafter\expand after\expandafter@gobble` after *A*, and seven `\expandafters` after *B*. But `\expandafter` can be completely avoided using a different approach, based on macro arguments:

```
\def\helpif#1#2{#1@firstoftwo}
\def\ucabs#1{%
  \if\noexpand#1.\helpif\fi@secondoftwo{
  {\if\noexpand#1a\helpif\fi@secondoftwo{A}
  {\if\noexpand#1b\helpif\fi@secondoftwo{B}
  {\noexpand#1\ucabs}}}%
}\message{\ucabs cbbas} % -> cBbas
```

It is not possible to move `\fi` into the definition of `\helpif`, because then \TeX won't see that particular `\fi` when it is skipping the whole `\if...` `\helpif` construction. With a small rearrangement we can get rid of `@secondoftwo`:


```

@gobble{\iftrue\iftrue}
  % \def\helpjf... contains 2*\fi
\def\helpjf#1\fi#1\expandafter\@firstoftwo
  \else\expandafter\@secondoftwo\fi}
\def\ucabj#1{%
  \helpjf\if\noexpand#1.\fi{}
  {\helpjf\if\noexpand#1a\fi{A}
  {\helpjf\if\noexpand#1b\fi{B}
  {\noexpand#1\ucabs}}}%
}\message{\ucabj cdabs} % -> cdAbs

```

The line containing `@gobble` above is needed so that `\def\helpjf` can be put inside `\iffalse ... \fi`, and \TeX 's `\fi` count won't decrease when seeing the two `\fi` tokens in the definition of `\helpjf`.

3.8 Appending Tokens to a Macro

\TeX doesn't provide primitives for modifying the expansion text of a macro, but it is possible to define a new macro with the contents of the old one and some additional tokens. For example, `\expandafter\def\expandafter\foo\expandafter{\foo\iffalse$}` appends the two tokens `\iffalse$` to the macro `\foo`. Three `\expandafter`s were used to make the old `\foo` only expanded once. None of the above tokens were expanded, fortunately. Any tokens can be appended this way, as long as they are nested with respect to braces. But care has to be taken when doubling `#s`:

```

\def\AppendTo#1#2{\expandafter\def\expandafter
  #1\expandafter{#1#2}}
\def\foo{} \AppendTo\foo{#}

```

yields the \TeX error "Illegal parameter number in definition of `\foo`". This can be solved by using token list registers, which double their hashmarks when expanded in an `\edef`:

```

\def\AppendTo#1#2{\begingroup
  \expandafter\toks\expandafter0%
  \expandafter{\foo#2}%
  \global\edef#1{\the\toks0}\endgroup}
\def\foo{} \AppendTo\foo{#x}
\show\foo % \foo=macro:-> ##x

```

Note that when a macro is defined, `#s` have to be properly formulated, and it is `\def` which eats half of `#s` (and converts `#1` etc. to special, inaccessible tokens), but `\edef` doesn't convert `#6` to special tokens if it comes from a token list register. The disadvantage of the second definition of `\AppendTo` is that it must be `\global`. (There is a much longer solution that manually doubles the `#s`.) The `\addto` command of Babel and the `\vref@addto` command of `varioref.sty` are `\global`, similar to this solution.

3.9 Processing Arbitrary Package Options

\LaTeX packages can use the standard commands `\DeclareOption`, `\ExecuteOption` and `\ProcessOptions` to access package options passed to them, and these commands work fine with a fixed set of options. The `\DeclareOption*` command can be used to declare arbitrary options:

```

%\DeclareOption{10pt}{\typeout{got ten-pt}}%(1)
\DeclareOption*{\typeout{got=(\CurrentOption)}}
\ProcessOptions % in file foo.sty

```

Two lines will be printed when `foo.sty` is loaded as `\usepackage[,foo=bar, no,]{foo}`. These are `got=(foo=bar)` and `got=no`. The optional argument of `\usepackage` may contain spaces and/or a single newline around commas and at the ends. Class options are passed to `\DeclareOption*`, so when `\documentclass[10pt]{article}` is active, `got=(10pt)` will not appear, but when line (1) is uncommented, `got ten-pt` will appear.

There is an alternative, low-level way for accessing all the options at once:

```

\AtEndOfPackage{\let\@unprocessedoptions
  \relax}% prevent warning
\typeout{\csname opt@\@currname.\@currentx
  \endcsname}

```

This prints the full option list with extra spaces and newlines removed, but commas, including superfluous ones, are kept intact.

4 Beyond the Current magyar.ldf

4.1 Other Hungarian Typesetting Software

Although `magyar.ldf` contains most of the functionality needed for following Hungarian typographic traditions, other utilities and packages can help in typesetting Hungarian texts. Most of this software, including `magyar.ldf`, is going to be available under the name `Magyar \LaTeX` from <http://www.math.bme.hu/latex/>.

`magyar.ldf` The new Hungarian module for Babel. Version 1.5 was written by Péter Szabó beginning in the autumn of 2003.

`huhyph.tex` or `huhyph3.tex` The old (version 3) Hungarian hyphenation `\patterns` for the T1 encoding, written by Gyula Mayer in 1998. Part of most \TeX distributions. See section 2.1.

`huhyphc.tex` The new version of the Hungarian hyphenation `\patterns` for the T1 encoding, written by Gyula Mayer in 2002 [4]. Part of most \TeX distributions. Hyphenates foreign compound words on the subword boundary, e.g. *szin-kron*. See section 2.1.

huhyphf.tex The new version of the Hungarian hyphenation `\patterns` for the T1 encoding, written by Gyula Mayer in 2002 [4]. Part of most TeX distributions. Hyphenates foreign compound words phonetically, e.g. *szink-ron*. See section 2.1.

ccs_extract.pl A Perl script that helps with hyphenation of words containing long double Hungarian consonants. It finds all occurrences of such words in the document and lets the user decide whether to insert, for each unique word, the `magyar.ldf` shorthands for `\discretionary` breaks for long double consonants. The program was written by Péter Szabó in 2003.

lafmtgen.pl An easy-to-use Perl script that can generate format files (`.fmt`) containing all hyphenation patterns required by the specified LaTeX document. It has some other features related to generating and installing format files, and is to be used with the Unix TeX distribution. It was written by Péter Szabó in 2003.

huplain.bst BibTeX style file for Hungarian bibliographies, based on `plain.bst`. It encourages sharing the same `.bib` database between Hungarian and English documents. It follows the (simple) convention that the style of a bibliography depends on the language of the document containing the entry, not the language of the entry itself. It was written by Péter Szabó in 2003.

husort.pl A drop-in replacement of `makeindex` that follows the Hungarian standards of index sorting and typesetting. It is implemented as a Perl script. It was written by Péter Szabó in 2003.

magyar.xdy Hungarian style file for the Xindy index processing program. Implements the Hungarian sorting order and a Hungarian typesetting style. The implemented sorting order does not follow Hungarian rules as strictly and elegantly as `husort.pl`. It was written by Péter Szabó in 2003.

CM-Super The EC fonts in Type 1 format in T1 and various other encodings. It is not part of MagyarLaTeX, but is available from CTAN. It is useful for converting Hungarian text to PDF, so the generated PDF file will contain the EC fonts in Type 1 format, and will be rendered quickly and nicely by Acrobat Reader.

MagyarISpell The Hungarian language database of the `lspell` spell checker for Unix. On Debian systems, it can be installed with the command `apt-get install ihungarian`.

`lspell` has a TeX mode, which skips control sequences and comments when checking

TeX source. (Unfortunately, the arguments of `\begin{tabular}` and many other non-textual elements of LaTeX documents are not skipped.) `lspell` can be used interactively, but this method is not comfortable, and incremental checking is not possible.

`lspell` also has an interprocess communication protocol, through which it can be integrated into text editors. For example, Emacs has built-in `lspell` support to mark incorrect words visually. OpenOffice, LyX, editors in KDE and newer versions of Vim can do the same. MagyarISpell works fine in these editors. It is not part of MagyarLaTeX, but it is freely available.

Note, however, that both the database and the stemmer of MagyarISpell is far from perfect, but among the Hungarian spell checkers only this one works inside `lspell`, so only this can be easily integrated into editors.

MSpell A commercial Hungarian spell checker with a no-cost Linux download, developed by Morphologic (a company in Hungary that produces linguistic software). Doesn't have an interactive mode, but can replace `lspell` in inter-process communication mode. A shell script is provided that replaces the `ispell` command, so MSpell can be integrated into text editors more easily. It is not part of MagyarLaTeX.

HunSpell The successor of MagyarISpell, but based on a different spell checking architecture. It understands Hungarian much better than MagyarISpell, but since it is not based on `lspell`, it is harder to integrate into text editors. For example, it is not available from the Emacs spell checking menu, even if it is installed. It is not part of MagyarLaTeX, but it is freely available.

4.2 Future Work

Some features are still missing from `magyar.ldf`:

- `letter.cls` is not customized properly, the left indentation of the nested list environments is also not customized;
- a macro to emit numbers with groups of three digits separated is missing;
- `layout.sty` and many other packages don't have Hungarian captions yet;
- the shorthand `'` is not disabled in math mode to give `nath.sty` a chance to typeset H_{symm} with `\mathcal{H}_{\text{'symm}}`;
- `\hunnewlabel` should store `table`, `figure` or `equation`, so `\refstruc` can insert it;
- new fonts and/or methods should be developed in place of `'tty`;

- Hungarian typography needs a baseline grid, which is almost impossible to enforce in L^AT_EX;
- some of the separation symbols proposed for after `\paragraph` are not available yet;
- section titles should not be larger than normal text;
- bold fonts should be substituted for bold extended fonts, whenever available—and never with an error or warning message; page numbers should be removed from blank pages;
- the width of `\parindent` should be computed based on `\textwidth`;
- the length of the last line of a paragraph should not be too near to the right margin, especially if `\parindent = 0`;
- `\vskips` above and below sections should be reduced;
- providing the commands `\H` and `\.` for OT1-encoded typewriter fonts;
- `\MathReal` should be extended with physical units;
- virtual fonts to support `\umlautlow` in T1 encoding;
- bold in `\begin{description}` should be `\emph`;
- allow specifying some compile-time options at run-time;
- `\textqq` should also work as an environment;
- `\told` should generate suffixes for month names.

Other features should be implemented outside `magyar.ldf`, as external programs. All programs in section 4.1 need improvement in one way or another.

4.3 Conclusion

An updated `magyar.ldf` which closely follows Hungarian typographical rules and works together with the most popular L^AT_EX packages without problems, has been awaited for many years. This new version is ready, as a single file longer than anything before, and it is filled with many advanced features.

Most of the features adapt L^AT_EX to Hungarian typographical rules, but some of them are bug fixes to various external packages, including design flaws and compatibility issues in **Babel** itself.

The implementation of some features clearly shows that T_EX macro programming is an obscure and ineffective way of solving some of the language-related problems. It is hoped that new versions of Ω , together with the new version of **Babel**, will provide a framework in which such problems can be addressed compactly and elegantly, without constant awareness of actual and possible compatibility glitches.

References

- [1] Gyöngyi Bujdosó and Ferenc Wettl. On the localization of T_EX in Hungary. *TUGboat*, 23(1):21–26, 2002.
- [2] David Kastrup. De ore leonis. Macro expansion for *virtuosi*. In *EuroBachoT_EX*, May 2002.
- [3] Donald E. Knuth. *The T_EXbook*. Addison-Wesley, 1984.
- [4] Gyula Mayer. The Hungarian hyphenation module of T_EX and L^AT_EX. Unpublished article in Hungarian, 10 July 2002.
- [5] Petr Sojka. Notes on compound word hyphenation in T_EX. *TUGboat*, 16(3):290–296, 1995.

Typesetting the Deseret Alphabet with L^AT_EX and METAFONT

Kenneth R. Beesley

Xerox Research Centre Europe

6, chemin de Maupertuis

38240 Meylan

France

Ken.Beesley@xrce.xerox.com

<http://www.xrce.xerox.com/people/beesley/>

Abstract

The Deseret Alphabet was an orthographical reform for English, promoted by the Church of Jesus Christ of Latter-day Saints (the Mormons) between about 1854 and 1875. An offshoot of the Pitman phonotypy reforms, the Deseret Alphabet is remembered mainly for its use of non-Roman glyphs. Though ultimately rejected, the Deseret Alphabet was used in four printed books, numerous newspaper articles, several unprinted book manuscripts, journals, meeting minutes, letters and even a gold coin, a tombstone and an early English-to-Hopi vocabulary. This paper reviews the history of the Deseret Alphabet, its Unicode implementation, fonts both metal and digital, and projects involving the typesetting of Deseret Alphabet texts.

1 Introduction

The Deseret Alphabet was an orthographical reform for English, promoted by the Church of Jesus Christ of Latter-day Saints (the Mormons) between about 1854 and 1875. While the Deseret Alphabet is usually remembered today as an oddity, a strange non-Roman alphabet that seemed doomed to failure, it was in fact used on and off for 20 years, leaving four printed books (including *The Book of Mormon*), numerous newspaper articles, several unprinted book manuscripts (including the entire *Bible*), journals, meeting minutes, letters and even a gold coin and a tombstone. There is also growing evidence that the Deseret Alphabet was experimentally used by some Mormon missionaries to transcribe words in Spanish, Shoshone, Hopi and other languages.

The Deseret Alphabet has been analyzed by a number of historians [19, 11, 20, 21, 4, 1, 22, 6] and justly criticized by typographers [21, 31], but what is often overlooked is the corpus of phonemically written documents, which are potentially interesting to both historians and linguists. Because few people, then or now, can be persuaded to learn the Alphabet, the majority of the documents have lain unread for 140 years. For example, in December of 2002, an “Indian Vocabulary” of almost 500 entries, written completely in the Deseret Alphabet, was finally identified as being English-to-Hopi, being perhaps the oldest written record of the Hopi language.

This paper will proceed with a short history of the Deseret Alphabet, putting it in the context of the Pitman phonotypy movement that inspired it from beginning to end;¹ special emphasis will be placed on the variants of the Alphabet used over the years, and on the cutting and casting of historical fonts. Then I will review some modern digital fonts and the implementation of the Deseret Alphabet in Unicode, showing how some honest mistakes were made and how the results are still awkward for encoding and typesetting some of the most interesting historical documents. Finally, I will show how I have used a combination of XML, L^AT_EX, the TIPA package and my own METAFONT-defined [16, 10] `desalph` font to typeset a critical edition of the English-to-Hopi vocabulary, and related documents, from 1859–60.

2 The Pitman Reform Context

2.1 The Pitman Reform Movements

To begin, it is impossible to understand the Deseret Alphabet without knowing a bit about two nineteenth-century orthographic reformers, Isaac Pitman (1813–1897) and his younger brother Benn (1822–1910). The Mormon experiments in orthographical reform, too often treated as isolated aberrations, were in fact influenced from beginning to

¹ Parts of this paper were first presented at the 22nd International Unicode Conference in San Jose, California, 11-13 September 2002 [6].

THE DESERET ALPHABET.

⊖	⊖	⊖	⊖	⊖	⊖
Long E.	A	AH	AW	O	OO
⊖	⊖	⊖	⊖	⊖	⊖
Short E.	A	AH	AW	O	OO
⊖	⊖	⊖	⊖	⊖	⊖
I	OW	WOO	YE	H	
7	⊖	⊖	⊖	⊖	⊖
F	B	T	D	CHE	G
⊖	⊖	⊖	⊖	⊖	⊖
K	GA	F	V	ETH	THE
⊖	⊖	⊖	⊖		
S	Z	USH	ZREE		
⊖	⊖	⊖	⊖	⊖	
R	L	M	N	NG	

WORDS OF ONE SYLLABLE.

⊖	⊖	⊖	⊖	⊖	⊖
EACH	AIM	AFT	ALL	OATH	OOZE
⊖	⊖	⊖	⊖	⊖	⊖
INK	EDGE	AM	ON	UP	FOOT
⊖	⊖	⊖	⊖	⊖	⊖
EYE	OWL	WOOD	YIELD	HANK	
⊖	⊖	⊖	⊖	⊖	⊖
PLOUGH	BOUGHT	TWIST	DROUGHT		
⊖	⊖	⊖	⊖	⊖	⊖
CHEAT	GRASS	CREEK	GAIN		
⊖	⊖	⊖	⊖	⊖	⊖
FACH	VERSE	MIRYATH	THOUGHT		
⊖	⊖	⊖	⊖	⊖	⊖
SALNT	ZEST	DISH	AZURE		
⊖	⊖	⊖	⊖	⊖	⊖
RIGHT	LAUGH	AGON	SING		

Figure 1: On 24 March 1854 the newly adopted Deseret Alphabet was first printed, probably using wooden type, and presented to the Board of Regents of the Deseret University. Although this rare flier is undated, it matches the 38-letter Alphabet as copied into the journal of Regent Hosea Stout on that date [30]. Utah State Historical Society.

end by the Pitman movements, at a time when many spelling reforms were being promoted.

2.1.1 Pitman Shorthand or Phonography

There have been hundreds of systems of stenography, commonly called shorthand, used for writing English; but Isaac Pitman's system, first published in his 1837 *Stenographic Sound-hand* and called "phonography",² was soon a huge success, spreading through the English-speaking world and eventually being adapted to some fifteen other languages. Modern versions of Pitman shorthand are still used in Britain, Canada, and in most of the cricket-playing countries; in the USA it was taught at least into the 1930s but was eventually overtaken by the Gregg system.

The main goal of any shorthand system is to allow a trained practitioner, called a "reporter" in the Pitman tradition, to record speech accurately at speed, including trial proceedings,³ parliamentary debates, political speeches, sermons, etc. Pitman's phonography, as the name implies, differs from most earlier systems in representing the distinctive sounds of English (what modern linguists call phonemes)

² In 1839 he wrote *Phonography, or Writing by Sound, being also a New and Improved System of Short Hand.*

³ In modern parlance we still have the term *court reporter*.

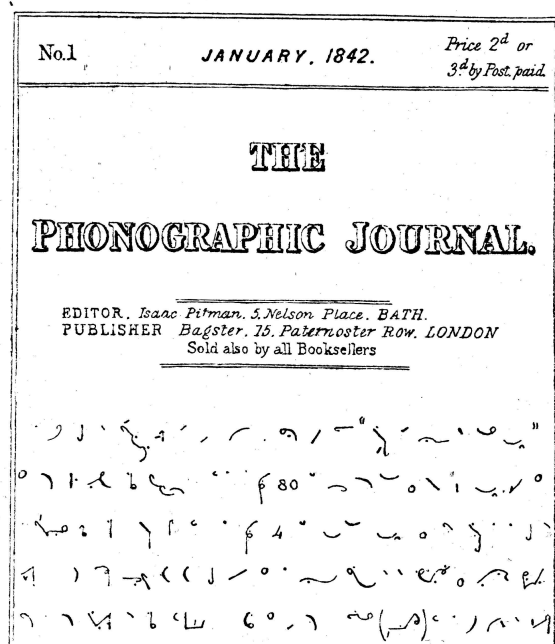


Figure 2: Early Pitman phonography.

rather than orthographical combinations. Simplicity and economy at writing time are crucial: consonants are reduced to straight lines and simple curves (see Figure 2). The “outline” of a word, typically just a string of consonants, is written as a single connected stroke, without lifting the pen. Voiced consonants are written as thick lines, their unvoiced counterparts as thin lines, which requires that a Pitman reporter use a nib pen or soft pencil. Vowels are written optionally as diacritical marks above and below the consonant strokes; one is struck by the similarities to Arabic orthography. In advanced styles, vowels are left out whenever possible, and special abbreviation signs are used for whole syllables, common words, and even whole phrases.

2.1.2 Pitman Phonotypy

Pitman became justifiably famous for his phonography. With help from several family members, he soon presided over a lecturing and publishing industry with a phenomenal output, including textbooks, dictionaries, correspondence courses, journals, and even books published in shorthand, including selections from Dickens, the tales of Sherlock Holmes, *Gulliver’s Travels*, *Paradise Lost*, and the entire *Bible*. But while phonography was clearly useful, and was both a social and financial success, Pitman’s biographers [25, 24, 2] make it clear that his real mission in life was not phonography but *phonotypy*,⁴ his philosophy and movement for reforming English orthography, the everyday script used in books, magazines, newspapers, personal correspondence, etc.

The first Pitman phonotypy alphabet for which type was cast was Alphabet No. 11, demonstrated proudly in *The Phonotypic Journal* of January 1844 (see Figure 3). Note that this 1844 alphabet is bicameral, sometimes characterized as an alphabet of capitals; that is, the uppercase and lowercase letters differ only in size. The letters are stylized, still mostly recognizable as Roman, but with numerous invented, borrowed or modified letters for pure vowels, diphthongs, and the consonants /θ/, /ð/, /ʃ/, /ʒ/, /tʃ/, /dʒ/ and /ŋ/.⁵

⁴ According to one of Pitman’s own early scripts, which indicates stress, he pronounced the word /foˈnɒtɪpi/.

⁵ To provide a faithful representation of original Pitman and Deseret Alphabet texts, I adopt a broad phonemic transcription that uses, as far as possible, a single International Phonetic Alphabet (IPA) letter for each English phoneme [12]. Thus the affricates ʧ and ʝ are transliterated as the rarely used IPA /tʃ/ and /dʒ/ letters, respectively, rather than the sequences /tʃ/ and /dʒ/ or even the tied forms /tʃ̥/ and /dʒ̥/. The diphthongs are shown in IPA as a combination of a nucleus and a superscript glide. The Deseret Alphabet, and the

The goals of general spelling reform, to create a new “book orthography”, are quite different from those of shorthand. While shorthand is intended for use by highly trained scribes, a book orthography is for all of us and should be easily learned and used. Where shorthand requires simplicity, abbreviation and swiftness of writing, varying with the reporter’s skill, a book orthography should aim for orthographical consistency, phonological completeness and ease of reading. Finally, a book orthography must lend itself to esthetic typography and easy typesetting; Pitman’s phonographic books, in contrast, had to be engraved and printed via the lithographic process.⁶

Pitman saw his popular phonography chiefly as the path leading to phonotypy, which was a much harder sell. His articles in the phonographic (shorthand) journals frequently pushed the spelling reform, and when invited to lecture on phonography, he reportedly managed to spend half the time talking about phonotypy. Throughout the rest of his life, Pitman proposed a long succession of alphabetic experiments, all of them Romanic, trying in vain to find a winning formula.

Pitman’s phonotypic publications include not only his phonotypic journals but dozens of books, including again the entire *Bible* (1850). But in the end, phonotypy never caught on, and the various phonotypic projects, including the constant cutting and casting of new type, were “from first to last a serious financial drain” [2]. In 1894, a few years before his death, Pitman was knighted by Queen Victoria for his life’s work in phonography, with no mention made of his beloved phonotypy.

Today Pitman phonotypy is almost completely forgotten, and it has not yet found a champion to sponsor its inclusion in Unicode. But Pitman was far from alone—by the 1880s, there were an estimated 50 different spelling reforms under consideration by the English Spelling Reform Association. This was the general nineteenth-century context in which the Deseret Alphabet was born; lots of people were trying to reform English orthography.

Pitman-Ellis 1847 alphabet which was its phonemic model, treat the /ju/ vowel in words like *mule* as a single diphthong phoneme; see Ladefoged [17] for a recent discussion and defense of this practice. Although in most English dialects the vowels in *mate* and *moat* are diphthongized, the Deseret Alphabet follows Pitman in treating them as the simple “long vowels” /e/ and /o/.

⁶ Starting in 1873, Pitman succeeded in printing phonography with movable type, but many custom outlines had to be engraved as the work progressed.

ADDRESS TO THE MEMBERS OF THE CORRESPONDING SOCIETY. 3

ADRE'S

TH AU MEMBURZ OV AU "FUNDGRAFIK KORESPONDIH SUSAL'ETI,"
AND AU SUBSKRABURZ TH AU FUNETIK FENT.

DIR FRENDRZ,—IT IZ WIA PLE'ZURABIL FELIHZ OV NU O'RDINURI KAND
DAA A ADRE'S W IN FUNDOTIPI, AND AUS OFUR W AU RIZU'LT OV AU FURST
EKSPERIMENT MED WIA AU FENT HWIE WR LIBERALITI HAZ ENEBILD
MI TH PROVA'D. TH W WIL FWCUR EJIZ LUK, AZ BHI, UNDER DIVA'N
PROVIDENS, AI INTRUDUSURZ OV A KOREKT MED OV RATH AND
PRINTIH : AI INSTRUKTURZ OV AU SIVILAZD WURLD IN AU TRU
PRINSIPILZ OV DAA ART HWIE IZ AU MENSPIIH OV SIVILIZE'SUN:
AI IMA'NSIPETERZ OV AI INFANT MAND FROM AU GOLIH CENZ OV AU
PREZENT SISTEM OV OROO'GRAFI : AND AI E'LVETURZ OV AU GRET MAS
OV MANKA'ND FROM AU LHEST DEPOS OV I'GNURANS AND SUPERSTI'SUN
TH AU PLEZURZ OV SAENS, AND AU DELA'TS OV VUREW.

Figure 3: In January 1844, Isaac Pitman proudly printed the first examples of his phonotypy. This Alphabet No. 11, and the five experimental variants that followed it, were bicameral, with uppercase and lowercase characters distinguished only by size.

2.2 The Mormons Discover the Pitman Movement

The Church of Jesus Christ of Latter-day Saints was founded in 1830 in upstate New York by Joseph Smith, a farm boy who claimed to have received a vision of God the Father and Jesus Christ, who commanded him to restore the true Church of Christ. He also claimed that he received from an angel a book, engraved on golden plates, which he translated as *The Book of Mormon*. His followers, who revered him as a prophet, grew rapidly in number, and soon, following the western movement and spurred by religious persecution, they migrated from New York, to Ohio, to Missouri and then to Illinois, where in 1839 they founded the city of Nauvoo on the Mississippi River.

Missionary work had started immediately, both at home and abroad, and in 1837, the same year that Pitman published his *Stenographic Sound-hand*, a certain George D. Watt was baptized as the first Mormon convert in England. Despite an unpromising childhood, which included time in a workhouse, young George had learned to read and write; and between the time of his baptism and his emigration to Nauvoo in 1842, he had also learned Pitman phonography. The arrival of Watt in Nauvoo revolutionized the reporting of Mormon meeting minutes, speeches and sermons. Other converts flowed into Nauvoo, so that by 1846 it had become, by some reports, the largest city in Illinois, with some 20,000 inhabitants.

But violence broke out between the Mormons and their “gentile” neighbors, and in 1844 Joseph Smith was assassinated by a mob. In 1845, even during the ensuing confusion and power struggles, Watt gave phonography classes; one notable student was Mormon Apostle Brigham Young. Watt was also President of the Phonographic Club of Nauvoo [1]. In addition to phonography, Watt was almost certainly aware of the new phonotypy being proposed by Pitman, and it is likely that he planted the idea of spelling reform in Brigham Young’s mind at this time.

In 1846, Watt was sent on a mission back to England. The majority of the Church regrouped behind Brigham Young, abandoned their city to the mobs, and crossed the Mississippi River to spend the bleak winter of 1846–47 at Winter Quarters, near modern Florence, Nebraska. From here Brigham Young wrote to Watt in April 1847:⁷

It is the wish of the council, that you procure 200 lbs of phonotype, or thereabouts, as you may find necessary, to print a small book for the benefit of the Saints and cause same to be forwarded to Winter Quarters before navigation closes, by some trusty brother on his return, so that we have the type to use next winter.

The “phonotype” referred to is the actual lead type used for Pitman phonotypy. The Saints, meaning

⁷ *The Latter-day Saints’ Millennial Star*, vol. 11, 1847, p. 8.

the members of the Church, were still in desperate times—600 would die from exposure and disease at Winter Quarters—and while there is no record that this type was ever delivered, it shows that the Mormons' first extant plans for spelling reform involved nothing more exotic than an off-the-shelf Pitman phonotypy alphabet.

298 PHONOTYPIC ALPHABET.

PHONOTYPIC ALPHABET.—No. 15.

VOWELS.				CONSONANTS.		
No.	Type.	Example of its sound.	Name in Phonotypes.	Type.	Example of its sound.	Name in Phonotypes.
1	J i	feet	i	P p	pay	pi
	I i	fit	it	B b	bay	bi
2	E e	mate	e	T t	toe	ti
	E e	met	et	D d	doe	di
3	A a	psalm	a	G g	chew	ge
	A a	Sam	at	J j	jew	je
4	Θ o	caught	o	K k	call	ke
	O o	cot	ot	G g	gall	ge
5	E c	cur	c	F f	few	ef
	U u	curry	ut	V v	view	ve
6	O o	bone	o	T t	thigh	it
	O o	fool	o	Ď đ	thy	đi
7	U u	full	ut	S s	seal	es
				Z z	zeal	ze
COMPOUND VOWELS.				Σ j	mesh	ij
Φ q	high	i		Ʒ ʒ	measure	ʒi
Φ q	hoy	q		L l	bail	el
U u	how	u		R r	bare	er
U u	hew	u		M m	sum	am
COALESCENTS.				N n	sun	en
Y y	yea	ye		U u	sung	ij
W w	way	we				
BREATHING.						
H h	hay	hc				

Figure 4: Alphabet No. 15 appeared in October 1844 and was the first of Pitman's "lowercase" or "small letter" alphabets, employing separate glyphs for uppercase and lowercase letters.

It is not known exactly which version of Pitman phonotypy Young had in mind; Pitman's alphabets went through no fewer than 15 variations between January 1844 and January 1847, and the isolated Mormons were likely out of date. In any case, Pitman's alphabets had by this time become more conventionally Roman. Alphabet No. 15 (see

Figure 4), presented in *The Phonotypic Journal* of October 1844,⁸ marked Pitman's abandonment of the bicameral "capital" alphabets, and his adoption of alphabets that had distinguished uppercase vs. lowercase glyphs, which he called "lowercase" or "small letter" alphabets.

The Mormons started leaving Winter Quarters as soon as the trails were passable, and the first party, including Brigham Young, arrived in the valley of the Great Salt Lake in July of 1847, founding Great Salt Lake City. Mormon colonists were soon sent throughout the mountain west. They called their new land Deseret, a word from *The Book of Mormon* meaning honey bee. In response to Mormon petitions to found a State of Deseret, Congress established instead a Territory of Utah, naming it after the local Ute Indians. In spite of this nominal rebuff, Brigham Young was appointed Governor, and the name Deseret would be applied to a newspaper, a bank, a university, numerous businesses and associations, and even a spelling-reform alphabet. The name Deseret, and the beehive symbol, remain common and largely secularized in Utah today.

3 The History of the Deseret Alphabet

3.1 Deliberations: 1850–1853

Education has always been a high priority for the Mormons, and on 13 March 1850 the Deseret University, now the University of Utah, was established under a Chancellor and Board of Regents that included the leading men of the new society. Actual teaching would not begin for several years, and the first task given to the Regents was to design and implement a spelling reform.

Although serious discussion of spelling reform began in 1850, I will jump ahead to 1853, when the Regency met regularly in a series of well-documented meetings leading to the adoption of the Deseret Alphabet. Throughout that year, the Regents presented to each other numerous candidate orthographies ranging from completely new alphabets, to Pitman shorthand, to minimal reforms that used only the traditional 26-letter Roman alphabet with standardized use of digraphs. The discussion was wide open, but by November of 1853, it was clear that the "1847 Alphabet" (see Figure 5), a 40-letter version backed jointly by Isaac Pitman and phonetician Alexander J. Ellis [15], was the recommended model. The 1847 Alphabet was presented to the Board in a surviving chart (see Figure 6) and the meeting minutes were even being delivered by

⁸ *The Phonotypic Journal*, vol. 3, no. 35, Oct. 1844.

THE PHONETIC ALPHABET.

The phonetic letters in the first column are pronounced like the italic letters in the words that follow. The last column contains the names of the letters.

<i>Long Vowels.</i>		<i>Explodents.</i>	
Ɛ	ε . . . ease ε	P	p . . . pole pe
Ɔ	a . . . age a	B	b . . . bowl be
Ȧ	q . . . alms q	T	t . . . toe te
Ɔ	ē . . . awning . . . ē	D	d . . . doe de
Ɔ	ō . . . ope o	Ɔ	g . . . cheer ga
W	u . . . ooze u	J	j . . . jeer ja
<i>Short Vowels.</i>		C	c . . . came ca
I	i . . . is it	G	g . . . game ga
E	e . . . egg et	<i>Continuants.</i>	
A	α . . . am at	F	f . . . fear ef
O	o . . . on ot	V	v . . . veer va
U	u . . . up ut	T	t . . . thigh it
W	u . . . sugar ut	Ɔ	ð . . . thy ðe
<i>Diphthongs.</i>		S	s . . . seal es
Ɔ	i . . . ice i	Z	z . . . zeal za
Ɔ	ē . . . oyster . . . ē	Σ	f . . . shall if
Ɔ	ɹ . . . ounce ɹ	Ɔ	ʒ . . . vision ʒe
Ɔ	u . . . use u	<i>Liquids.</i>	
<i>Coalescents.</i>		R	r . . . rare ur
Y	y . . . yea ya	L	l . . . lull el
W	w . . . way wa	<i>Nasals.</i>	
<i>Breathing.</i>		M	m . . . mum cm
H	h . . . hay hag	N	n . . . nun en
		W	ŋ . . . sing ij

(*) *Vocal*, as in *ab'l*, *siz'm*, *hev'n*, &c.

Figure 5: The 1847 Alphabet of Alexander J. Ellis and Isaac Pitman as it appeared in Pitman's 1850 *Bible*. This alphabet was the main phonemic model for the Deseret Alphabet in late 1853. The Board of Regents of the Deseret University almost adopted a slightly modified form of this alphabet, but they were persuaded, at the very last moment, to change to non-Roman glyphs. Compare the layout of this chart to that of the Deseret Alphabet charts in the books of 1868–69 (see Figure 17).

reporter George D. Watt in the longhand form of this alphabet.

Brigham Young, President of the Church of Jesus Christ of Latter-day Saints and Governor of the Territory of Utah, took a personal interest in the 1853 meetings, attending many and participating actively. On the 22nd and 23rd of November, he and the Regents adopted their own modified version of the 1847 Alphabet, with some of the glyphs

modified or switched, and names for the letters were adopted. A couple of Pitman letters were simply voted out, namely those for the diphthongs /^ji/ and /^ju/, which are exemplified with the words *oyster* and *use* in the 1847 chart. The result was a 38-letter alphabet, still very Pitmanesque and Romanic. For the second time—the first was in 1847—the Mormons were about to embark on a Pitman-based spelling reform.

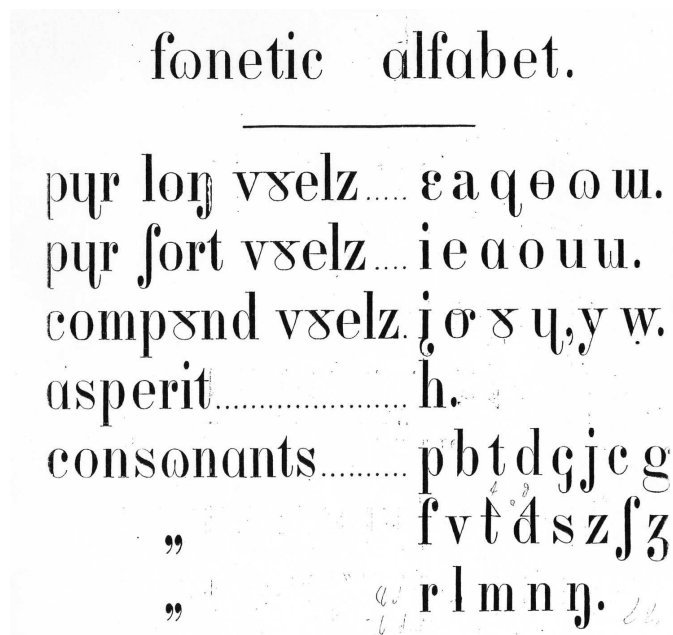


Figure 6: In November 1853, Parley P. Pratt presented “Pitman’s Alphabet in Small Letters” to the Board of Regents in the form of this chart. These are in fact just the lowercase letters of the famous 1847 Alphabet devised by Isaac Pitman and Alexander J. Ellis. More stable than Pitman’s other alphabets, it lasted several years and was used to print a short-lived newspaper called *The Phonetic News* (1849), the *Bible* (1850), and other books. LDS Church Archives.

However, all plans were turned upside-down by the sudden arrival of Willard Richards at the meeting of 29 November 1853. Richards, who was Second Counselor to Brigham Young, was gravely ill, had not attended the previous meetings, and was not up to date on the Board’s plans. But when he saw the Board’s new Romanic alphabet on the wall, he could not contain his disappointment. The following excerpts, shown here in equivalent IPA to give the flavor of George D. Watt’s original minutes, speak for themselves:

wi wɔnt e nju kəˈnd ðv ælfæbet, difəriŋ frɔm ði kɔmpaˈnd mɛs ðv stɑf ʌpɔn ðæt fit.... ðoz kæræktɛrz mɛ bi ɛmplɔɪd ɪn ɪmpɹuvɪŋ ði ŋ-ɡlɪʃ ɔrθɔgræfɪ, ðo æt ði sɛm tɑˈm, ɪt ɪz æz əˈj hæv sɑmtɑˈmz sɛd, ɪt sɪmz lɑˈk pɑtɪŋ nju wɑˈn ɪntu old bɔtlz.... əˈj æm ɪnklaɪnd tu θɪŋk hwɛn wi hæv rɪflɛktɛd hɔŋɛr wi ʃæl stɪl mɛk sɑm ædvæns ʌpɔn ðæt ælfæbet, ænd prhæps θro æwɛ ɔl kæræktɛrz ðæt bɛr mɑtʃ rɪzɛm-blɛns tu ði ŋɡlɪʃ kæræktɛrs, ænd ɪntrodju s æn ælfæbet ðæt ɪz ɔrɪʒɪmæl, so fɑr æz wi no, æn ælfæbet ɛntɑˈrɪli dɪfɛrɛnt frɔm ɛnɪ ælfæbet ɪn juː.⁹

⁹ “We want a new kind of alphabet, differing from the compound mess of stuff upon that sheet. ... Those charac-

Some objections were tentatively raised. It was pointed out that the key committee had been instructed to keep as many of the traditional Roman letters as possible, and that Brigham Young himself had approved the alphabet and had already discussed ordering 200 pounds of type for it. Richards then attenuated his criticism a bit, but renewed his call for a complete redesign, waxing rhetorical:

hwɔt hæv ju gɛnd bɑˈj ði ælfæbet ɔn ðæt kɑrd əˈj æsk ju. ʃo mi wʌn əˈtɛm, kæn ju pɔɪnt əˈwɪt ði fɛrst ædvæntɛdʒ ðæt ju hæv gɛnd ɔvɛr ði old wʌn? ... hwɔt hæv ju gɛnd, ju hæv ði sɛm old ælfæbet ɔvɛr æŋɡɛn, ɔnlɪ ə ʃu ædɪf-næl mɑrks, ænd ðɛ ɔnlɪ mɪstɪfɑˈj ɪt mɔr, ænd mɔr.¹⁰

ters may be employed in improving the English orthography, though at the same time, it is as I have sometimes said, it seems like putting new wine into old bottles. ... I am inclined to think when we have reflected longer we shall still make some advance upon that alphabet, and perhaps throw away all characters that bear much resemblance to the English characters, and introduce an alphabet that is original, so far as we know, an alphabet entirely different from any alphabet in use.”

¹⁰ “What have you gained by the alphabet on that card I ask you. Show me one item, can you point out the first advantage that you have gained over the old one? ... What have you gained, you have the same old alphabet over again,

Richards believed fervently that the old Roman letters varied too much in their values, that no one would ever agree on their fixed use, and that keeping them would just be a hindrance; a successful, lasting reform would require starting with a clean slate. He also argued for economy in writing time, paper and ink. These arguments anticipated those advanced by George Bernard Shaw in the 20th century to support the creation of what is now known as the Shaw or Shavian Alphabet [28, 18].¹¹

Brigham Young and the Board of Regents were persuaded, the Board's modified Pitman alphabet was defenestrated, and the first version of a new non-Roman alphabet was adopted 22 December 1853, with 38 original glyphs devised by George D. Watt and perhaps also by a lesser-known figure named John Vance. The Deseret Alphabet was born.

3.2 Early Deseret Alphabet: 1854–1855

In Salt Lake City, the *Deseret News* announced the Alphabet to its readers 19 January 1854:

The Board of Regents, in company with the Governor and heads of departments, have adopted a new alphabet, consisting of 38 characters. The Board have held frequent sittings this winter, with the sanguine hope of simplifying the English language, and especially its Orthography. After many fruitless attempts to render the common alphabet of the day subservient to their purpose, they found it expedient to invent an entirely new and original set of characters.

These characters are much more simple in their structure than the usual alphabetical characters; every superfluous mark supposable, is wholly excluded from them.

The written and printed hand are substantially merged into one.

Type of some kind, almost certainly wooden,¹² was soon prepared in Salt Lake City, and on 24 March 1854 a four-page folded leaflet with a chart of the Deseret Alphabet was presented to the Board (see Figure 1). In this early 1854 version of the Alphabet, we find 38 letters, the canonical glyphs being drawn with a broad pen, with thick emphasis on the downstrokes, and light upstrokes and flourishes. The short-vowel glyphs are represented smaller than the others.

George D. Watt was the principal architect of the Deseret Alphabet and, judging by surviving documents, only a few additional marks, and they only mystify it more, and more.”

¹¹ <http://www.shavian.org/>

¹² *Deseret News*, 15 August 1855.

uments, was also the first serious user. Watt was a Pitman stenographer, and the early documents (see Figure 7) are written in a distinctly stenographic style.¹³ Watt drew the outline of each word cursorily, without lifting the pen. Short vowels, shown smaller than the other glyphs in the chart, were incorporated into the linking strokes between the consonants; thus vowels were usually written on upstrokes, which explains their canonical thin strokes and shorter statures in the first chart. The writer had to go back and cross the † vowels after finishing the outline; and often short vowels were simply left out.

The demands of cursive writing seem to have influenced the design of several of the letters. In particular, the fussy little loops on the θ (/d/), δ (/s/), \mathfrak{G} (/g/), \mathfrak{O} (/o/) and \mathfrak{D} (/a^w/) were used to link these letters with their neighbors. Watt also combined consonants together with virtuosity, “amalgamating” them together to save space, but at the expense of legibility. Another lamentable characteristic of the early style was the inconsistent use of letters, sometimes to represent their phonemic value and sometimes to represent their conventional name. Thus Watt writes *people* as the equivalent of /ppl/, expecting the reader to pronounce the first p-letter as /pi/, that being the letter's conventional name when the alphabet is recited. Similarly, Watt can spell *being* as the equivalent of just /bɨ/, the letters having names pronounced /bi/ and /ɨj/, respectively. While probably seen by shorthand writers as a clever way to abbreviate and speed their writing, the confusion of letter names and letter values is a mistake in any book orthography.

Like Isaac Pitman, the Mormons could not resist experimenting with their new alphabet, changing both the inventory of letters and the glyphs. The 1854 alphabet was almost immediately modified, substituting new glyphs for /i/ and /a^w/ and adding two new letters for the diphthongs /ɔ^j/ and /i^ju/, making a 40-letter alphabet as printed in the 1855 *Deseret Almanac* of W.W. Phelps. This chart was almost certainly the one copied by Rémy and Brenchley [27] who visited Salt Lake City in 1855 (see Figure 8).¹⁴

Watt apparently believed that the same basic alphabet could serve for both stenography and everyday orthography, or as the *Deseret News*, cited above, put it, “The written and printed hand are

¹³ James Henry Martineau was another early cursive writer.

¹⁴ For yet another chart of this version of the Alphabet, see Benn Pitman's *The Phonographic Magazine*, 1856, pp. 102–103.

13⁴ uoaa soaa qm8
5 7. 7. uoaa3 qaa 6⁴ 1854.

uoyy8 n7 78077 70 aa 4m7 8 841a m 8m3 147.
 auoy 7uoo uue oaa 70 8 07 70 77 7uu auoy 7u7u⁷ 7u7a.
 "8 7u7 7u7 oaa 14 3 7u7 6 auu"; . . . 7u7 8m, 7u7 7u
 aa auoy 8707.
 8 7u7 7u7 8 1477 7m 7u7 7u7 7u7 7u7 7u7.
 auoy 7u7 7u7 7u7 7u7 7u7 7u7 7u7 7u7 7u7 7u7.
 8 1, 2, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17,
 u7 7u7, 7u7 7u7 7u7 7u7 7u7 7u7 7u7 (7u7) 7u7 7u7 7u7.
 auoy 7u7 7u7 7u7 7u7 7u7 7u7 7u7 7u7 7u7 7u7,
 7u7 7u7 7u7 7u7 7u7 7u7 7u7 7u7 7u7 7u7.
 7 7u7 8 7u7 7u7 8 7u7 7u7 7u7 7u7 7u7 7u7 7u7 7u7 7u7,
 7u7 7u7 7u7 7u7 7u7 7u7 7u7 7u7 7u7 7u7 7u7 7u7 7u7 7u7
 7u7 7u7. 7u7 7u7 7u7 7u7 7u7 7u7 7u7 7u7 7u7 7u7 7u7 7u7 7u7

Figure 7: Extract from the minutes of a Bishops’ meeting, 6 June 1854, concerning the support of the poor. These minutes, written in a cursive, stenographic style, were prepared by George D. Watt and addressed directly to Brigham Young. LDS Church Archives.

substantially merged into one.” This was in fact an early goal of phonotypy, but it was soon abandoned by Pitman as impractical [15]. The retention of this old idea contributed to making the Deseret Alphabet an esthetic and typographical failure.

One of the fundamental design problems in the Alphabet was the elimination of ascenders and descenders. This was done in a well-intentioned attempt to make the type last longer—type wears out during use, and the ascenders and descenders wear out first—but the lamentable result was that all typeset words have a roughly rectangular shape, and lines of Deseret printing become very monotonous. Some of the glyphs, in particular ϑ and \mathcal{Q} , are overly complicated; and in practice writers often confused the pairs \mathcal{Q} vs. \mathcal{Q} and \mathcal{Q} vs. \mathcal{D} . These fundamental design problems need to be distinguished from the font-design problems, which will be discussed below.

3.3 The 1857 St. Louis Font

The reform was moving a bit slowly. On 4 February 1856 the Regents appointed George D. Watt, Wilford Woodruff, and Samuel W. Richards to prepare manuscripts and arrange for the printing of books. The journals of Richards and Woodruff show that they went at it hammer and tongs, working on elementary readers and a catechism intended for teach-

ing religious principles to children. The next step was to get a font made.

There are references to an attempt, as early as 1855, to cut Deseret Alphabet punches right in Utah, by a “Brother Sabins”,¹⁵ but there is as yet no evidence that this project succeeded. In 1857, Erastus Snow was sent to St. Louis to procure type, engaging the services of Ladew & Peer, which was the only foundry there at the time [31]. But Snow abandoned the type and hurried back to Utah when he discovered that President Buchanan had dispatched General Albert Sydney Johnston to Utah with 2500 troops from Fort Leavenworth, Kansas, to put down a reported Mormon rebellion and install a new non-Mormon governor. The news of “Johnston’s Army” reached Salt Lake City 24 July 1857, when the alleged rebels were gathered for a picnic in a local canyon to celebrate the tenth anniversary of their arrival in Utah. In the ensuing panic, Salt Lake City and the other northern settlements were abandoned, and 30,000 people packed up their wagons and moved at least 45 miles south to Provo. The territorial government, including records and the printing press, were moved all the way to Fillmore in central Utah. While this bizarre and costly fiasco,

¹⁵ *The Latter-day Saints’ Millennium Star*, 10 November 1855. The reference is probably to John Sabin (not Sabins), who was a general mechanic and machinist.



Figure 8: Rémy and Brenchley almost certainly copied this chart from an almost identical one in W.W. Phelps’ *Deseret Almanac* of 1855. With the addition of letters for /ɔ^j/ and /j^u/, this 40-letter version of the Deseret Alphabet had the same phonemic inventory as the Pitman-Ellis 1847 Alphabet.

often called the Utah War or Buchanan’s Blunder, was eventually resolved peacefully, it was another setback to the Deseret Alphabet movement.

By late 1858, the Utah War was over, the St. Louis type had arrived in Salt Lake City, and work recommenced. It is very likely that only the punches and matrices were shipped to Utah,¹⁶ and that the Mormons did the actual type casting themselves. The children’s texts prepared by the committee of Woodruff, Richards and Watt had been lost; unfazed, Brigham Young told Woodruff to “take hold with Geo. D. Watt and get up some more”.¹⁷ The first use of the new type was to print a business card for George A. Smith, the Church Historian. The stage was now set for the revival of the Deseret Alphabet reform in 1859–60.

¹⁶ *Deseret News*, 16 February 1859.

¹⁷ *Journal History*, 20 November 1858. The journal of Wilford Woodruff for 22 November 1858 indicates that the manuscripts were soon found.

3.4 The Revival of 1859–1860

3.4.1 Sample Articles Printed in the *Deseret News*

The period of 1859–60 was a busy and productive one for the Deseret Alphabet. The type was finally available, and on 16 February 1859 the *Deseret News* printed a sample text from the Fifth Chapter of Matthew, the Sermon on the Mount. Similar practice texts, almost all of them scriptural, appeared almost every week to May 1860. Despite this progress, everyone involved was extremely disappointed with the St. Louis font, which was crudely cut and ugly by any standards. Young felt that the poor type did as much as anything to hold back the reform.

The 1859 Alphabet as printed in the *Deseret News* (see Figure 9) had reverted to 38 letters, lacking dedicated letters for the diphthongs /ɔ^j/ and /j^u/, which had to be printed with digraphs; but

		DESERET ALPHABET.			
Long.	Short.	Ƴ	h	L	eth
ə	e	7	p	ʒ	the
3	a	8	b	8	s
ə	ah	7	t	6	z
ə	au	8	d	D	esh
O	o	c	che	s	zhe
0	oo	9	g	4	ur
h	i	Q	k	l	l
ə	ow	Q	ga	7	m
U	woo	f	f	4	n
Ƴ	ye	8	v	U	eng

890 w6 704774, 739 228.

hə ʒ 83 ʔh79 Ƴ0, ʒ77 Ƴə 7ʔ87
 4+777, 748 8+ 8047 7074: Ƴw4
 ʒ+ 87+477 87L, 7 Ƴə 74 4w7
 8047 7074, Ƴə 074w7 747447
 ʒ+ 0+7077 76 Ƴ764; ʒ74Ƴ04 077
 743 8+ 8777768 ʔh79 4+7747748,
 ʒ77 Ƴə 73 8+ 7w07 Ƴ4w7 Ƴ04
 8+46, ʒ77 Ƴə 73 Ƴ76 Ƴ3L w4 ʒ+
 777 76 0w8, Ƴ0 7307L 7U3 ʒ+
 8+46 w6 ʒ+ 77478, Ƴ0 76 7777 79
 838 748 79 07746 Ƴ4w7 07 ʔh-
 4778778; Ƴ3, ʒ 83 ʔh79 Ƴ0, 077
 748 Ƴ4 4w7, 748 73 7878 7647
 8+4, Ƴ77c 06777 87L 8+877 Ƴ0,

Figure 9: The *Deseret News* printed sample articles in the Deseret Alphabet in 1859–60, and again in 1864, using the crude St. Louis type of 1857. This article, of which only a portion is shown here, appeared in the issue of 30 November 1864, vol. XIV, no. 9, which also included reports of the fall of Atlanta, Georgia to General Sherman during the American Civil War.

the *Deseret News* apologized for the lack of a /ju/ letter and promised a correction as soon as a new punch could be cut.¹⁸

In 2002 I found the punches for the 1857 St. Louis font in the LDS Church Archives (see Figure 10). There proved to be only 36 punches in

¹⁸ The *Deseret News* also promised a new letter for the vowel in *air*, which was a highly suspect distinction made in some Pitman alphabets.

each of three sizes, but investigation showed that they were originally intended to support a 40-letter version of the Alphabet. The trick was the double use of four of the punches, rotating them 180 degrees to strike a second matrix. Thus the punch for 7 also served to strike the matrix for ʒ; the punch for 7 also served for L; and similarly for the pairs ʒ-Ƴ and 7-7. The sets include punches for the /ɔʃ/ and /ju/ diphthongs, being 0 and 0, respectively, but these glyphs had apparently fallen out of favor by 1859 and were not used in the *Deseret News*.



Figure 10: Some smoke proofs of the 1857 St. Louis punches, found in 2002 in the LDS Church Archives. The 0 and 0 glyphs, representing /ɔʃ/ and /ju/, respectively, were not used when the *Deseret News* finally started printing sample articles with the type in 1859.

3.4.2 Handwritten Deseret Alphabet in 1859–60

Brigham Young directed his clerks to use the Alphabet, and the history or biography of Brigham Young was kept in Deseret Alphabet at this time. Another surviving text from this period is the financial “Ledger C”, now held at Utah State University (see Figure 12). This ledger was probably kept by clerk T.W. Ellerbeck who later wrote [19], “During one whole year the ledger accounts of President Young were kept by me in those characters, exclusively, except that the figures of the old style were used, not having been changed.”

The Ledger C alphabet has 39 letters, including the glyph ʔ for /ju/ but using a digraph for /ɔʃ/. The Ledger abandons the Alphabet in May of 1860, at the same time that the *Deseret News* stopped printing sample articles, and the Deseret text was at some point given interlinear glosses in standard orthography.

ᐱᐱᐱᐱ.

ᐱᐱᐱ 3, ᐱᐱᐱᐱ ᐱᐱᐱᐱ, ᐱ, ᐱ, ᐱᐱ ᐱ ᐱᐱᐱ ᐱ ᐱᐱᐱᐱᐱ

ᐱᐱᐱ 8ᐱᐱ8 ᐱᐱᐱ.

“ 8, 6ᐱ ᐱᐱᐱ ᐱᐱᐱ ᐱᐱᐱᐱ ᐱᐱᐱ ᐱᐱᐱᐱ.

“ 8, 10ᐱ “ “ ᐱᐱᐱ “ ᐱᐱᐱ.

“ 8, 12ᐱ “ “ ᐱᐱᐱ “ ᐱᐱᐱ.

“ 9, 11ᐱ “ “ ᐱᐱᐱ “ ᐱᐱᐱ.

“ 18, 3ᐱᐱ “ “ ᐱᐱᐱ “ ᐱᐱᐱ.

Figure 11: A portion of the errata sheet, printed in Utah using the St. Louis type of 1857, for *The Deseret First Book* of 1868. A much better font was cut for printing the readers (see Figure 15), but it was left in the care of Russell Bros. in New York City.

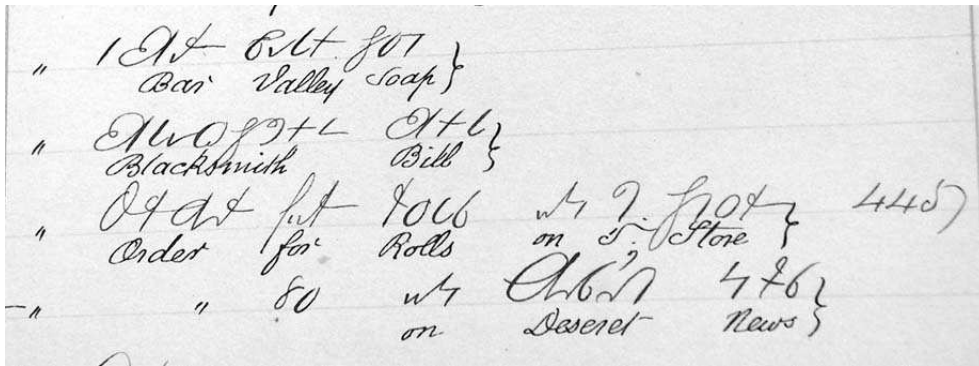


Figure 12: *Ledger C* of 1859-60, probably kept by T.W. Ellerbeck, illustrates an idiosyncratic 39-letter version of the Deseret Alphabet. There are still cursive links and “amalgamations” in this text, but far fewer than in George D. Watt’s early texts of 1854–55. Interlinear glosses in traditional orthography were added later. Utah State University.

My own favorite document from this era is the Deseret Alphabet journal of Thales Hastings Haskell [29], kept from 4 October 1859 to the end of that year while he and Marion Jackson Shelton¹⁹ were serving as Mormon missionaries to the Hopi.²⁰ They were staying in the Third-Mesa village of Orayvi (also spelled Oribe, Oraibi, etc.), now celebrated as the oldest continuously inhabited village in North America. Haskell used a 40-letter version of the alphabet, like the contemporary *Deseret News* version, but adding ʔ for /^ju/ and, idiosyncratically, 0 for /ɔ^j/. The original manuscript is faint and fragile; the following is a sample in typeset Deseret Alphabet and equivalent IPA:

¹⁹ Shelton also kept a journal in 1858–59, in a mix of standard orthography, Pitman shorthand, and some Deseret Alphabet. LDS Church Archives.

²⁰ The original journal is in Special Collections at the Brigham Young University Library. At some unknown time after the mission, Haskell himself transcribed the Deseret text into standard orthography, and this transcription was edited and published by Juanita Brooks in 1944 [9].

ᐱᐱᐱ ᐱᐱ ᐱᐱᐱ ᐱᐱᐱᐱᐱ ᐱᐱᐱ ᐱᐱᐱᐱᐱ ᐱᐱᐱᐱ ᐱᐱᐱᐱ ᐱᐱ ᐱᐱᐱᐱ
ᐱᐱᐱᐱ ᐱᐱ ᐱᐱᐱ ᐱᐱᐱ ᐱᐱᐱ ᐱᐱᐱ ᐱᐱᐱᐱᐱ

gɔt ʌp tɔk bɛkʌst [sic] ænd stɑtɪd mɪɑn
wɛnt æhɛd tu ɔrɑːb vɪlɪʒ tu tɛl ðɛm ðæt wi
wɛr kʌmɪŋ

In standard orthography, this reads, “Got up, took breakfast and started [sic] Indian went ahead to Oribe village to tell them that we were coming.” The missing *r* in *breakfast* is just an isolated error, but the spelling of /statid/ for *started* is characteristic; Haskell was from North New Salem, Franklin County, Massachusetts, and he dropped his *rs* after the /a/ vowel [5]. Other writers similarly leave clues to their accents in the phonemically written texts.

Marion J. Shelton was a typical 40-letter Deseret writer from this period, using the more or less standard glyphs 0 for /ɔ^j/ and ʔ for /^ju/, as in the

3.5 The 1860s and the Printed Books

Most of the enthusiasm for the Deseret Alphabet collapsed in 1860, and by 1862 it was dead, except in the determined mind of Brigham Young. When Superintendent of Common Schools Robert L. Campbell presented Brigham Young with a manuscript of a “first Reader” in standard orthography, Young rejected it emphatically, insisting that “he would not consent to have his type, ink or paper used to print such trash”.²⁴

Æ ALFABET

Ov ðs Kapitalz and Smol Letterz.

E	e	H	j	F	f
E	e	G	o	V	v
A	a	S	s	H	t
A	a	U	u	A	d
A	q	Y	y	S	s
O	o	W	w	Z	z
O	o	H	h	Σ	j
O	o	H	h	Σ	j
I	i	P	p	Σ	j
E	e	B	b	L	l
A	a	T	t	R	r
A	a	D	d	R	r
O	o	G	g	M	m
O	o	J	j	N	n
U	u	K	k	N	n
U	u	G	g	W	g

Figure 13: The 1855 Benn Pitman or American Pitman Alphabet. In 1852, Benn Pitman carried the Pitman phonography and phonotypy movement to the United States, setting up The Phonographic Institute in Cincinnati in 1853. Whereas Isaac Pitman was an incurable tinkerer, constantly modifying his alphabets, brother Benn recognized the virtues of stability.

In 1864, the Regents considered adopting the phonotypy of Benn Pitman, the brother of Isaac who had established his own Phonographic Institute in Cincinnati in 1853, but the ultimate response was a recommitment to the Deseret Alphabet; sample Alphabet articles reappeared defiantly in the *Deseret*

²⁴ Journal History, 22 May 1862.

News 11 May 1864 and continued to the end of the year.

There were in fact several attempts during the 1860s to abandon the Deseret Alphabet. In December of 1867,²⁵ the Board of Regents, with Brigham Young, resolved unanimously to adopt “the phonetic characters employed by Ben [sic] Pitman of Cincinnati, for printing purposes, thereby gaining the advantage of the books already printed in those phonetic characters.” However, on 3 February 1868,²⁶ the Board once again did an about-face, recommitted to the Deseret Alphabet and started the serious and expensive work of getting books prepared for publication. Apostle Orson Pratt was hired to transcribe *The Deseret First Book* and *The Deseret Second Book* into the Deseret Alphabet.

After the disappointment with the crude St. Louis type, the Regents in 1868 sent their agent D.O. Calder to New York to get better workmanship. Calder engaged the firm of Russell Bros.,²⁷ which cut and cast an English (14-point) font for the project. The new school books (see Figures 14 and 15) were delivered to Salt Lake City in late 1868, at which time Orson Pratt had already turned his dogged energy to the transcription of *The Book of Mormon*.

In 1869, Pratt was sent as the agent to New York, to supervise the printing of *The Book of Mormon*. He too chose Russell Bros. and had a font of Long Primer (10-point) type cut and cast for the body of the text.²⁸ The bicameral nature of the Deseret Alphabet allowed him to save some money by using the lowercase letters of the existing English (14-point) font as the uppercase letters of the Long Primer (10-point) font. Pratt also had fonts prepared in the Great Primer (18-point) and Double English (28-point) sizes to serve in headings and titles. Not surprisingly, Pratt complained that the three unlucky compositors assigned to the project were making “a great abundance of mistakes in setting type”, and he had to give the proofs four good readings and supervise many corrections before the pages could be stereotyped.²⁹

²⁵ *Deseret News*, 19 December 1867.

²⁶ *Deseret News*, 3 February 1868.

²⁷ Russell’s American Steam Printing House, located at 28, 30 and 32 Centre Street, New York City, Joseph and Theodore Russell, Props.

²⁸ Small Pica (11-point) type was also considered and, unfortunately, rejected. With the inherent design problems of the Deseret Alphabet, the Long Primer type is too small for comfortable reading.

²⁹ Orson Pratt to Robert L. Campbell, 12 June 1869, Deseret Alphabet Printing Files 1869, LDS Church Archives.



Figure 14: In 1868, *The Deseret First Book*, shown here, and *The Deseret Second Book* were printed by Russell Bros. of New York and shipped to the Territory of Utah. The print run for each book was 10,000 copies.

The Book of Mormon (see Figure 16) was published in two formats. *The Book of Mormon Part I*, intended to serve as an advanced reader, consisted of *The First Book of Nephi*, *The Second Book of Nephi*, *The Book of Jacob*, *The Book of Enos*, *The Book of Jarom*, *The Book of Omni* and *The Words of Mormon*.³⁰ The entire *Book of Mormon* was also printed on better paper, and was more expensively bound.

³⁰ *The Book of Mormon Part I* is usually known, inaccurately, among used-book dealers as “The First Book of Nephi”. The Regents’ plan was eventually to offer the whole book in three parts, printing Parts II and III with proceeds from the sale of the first four books.

Receipts from 1868 and 1869³¹ show that the punches, matrices, type and other printing paraphernalia remained the property of the Board of Regents of the Deseret University, but they were left in the care of Russell Bros. in expectation of future work, which in fact never materialized. Although a large collection of nineteenth-century punches survives at Columbia University in New York City, attempts to locate the Russell Bros. Deseret Alphabet punches have so far been unsuccessful.

³¹ Deseret Alphabet Printing Files 1868 and 1869, LDS Church Archives.

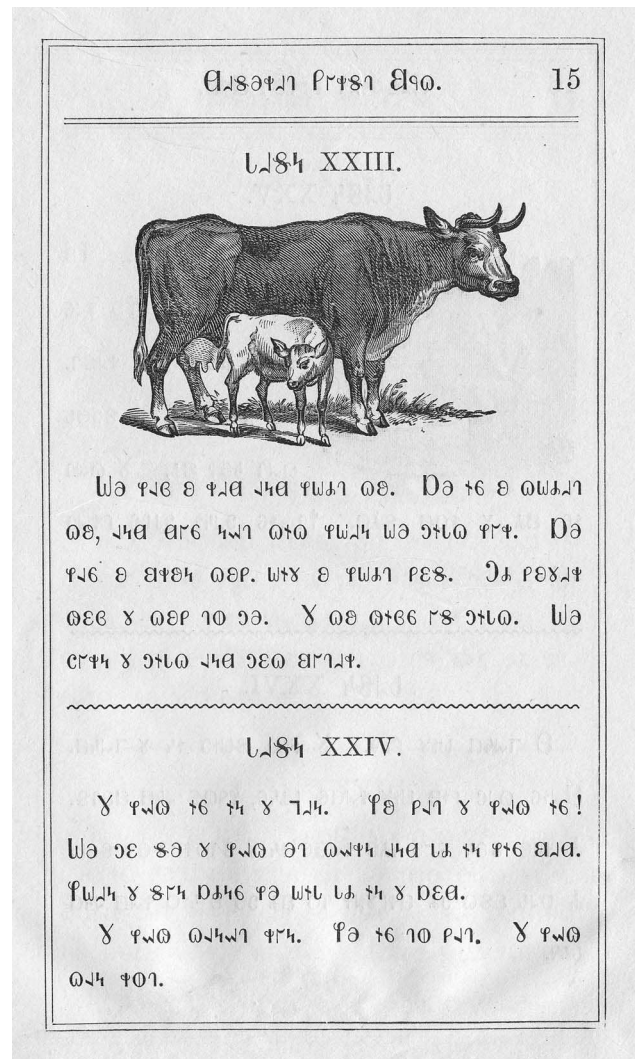


Figure 15: A page from *The Deseret First Book*.

3.6 The 1868–69 “Book” Alphabet and Fonts

After the disappointing debut of the St. Louis type, used reluctantly to print sample articles in the *Deseret News* in 1859–60 and 1864, Brigham Young had vowed to go to England the next time to get better workmanship.³² But in fact in 1868 and 1869 the Mormons went only as far as New York City, engaging Russell Bros. to cut new punches, strike matrices, cast type, typeset and print the books.

This time they did get professional workmanship, but the resulting book font is still somewhat bizarre, partly because of the inherent awkwardness of the basic shapes, and partly because of choices in font design that now seem old-fashioned. A look

at the book font (see Figures 15, 16 and 17) shows that the glyphs, compared to the earlier charts, have been Bodonified: made rigidly vertical, symmetrical wherever possible, and with extreme contrasts of thick and thin. Thom Hinckley, an expert typographer and printer (personal communication), has pointed out that the extreme thins of the font reveal the punch cutter as a master; at the same time, these thin lines would have caused the type to wear quickly, which was one of the very problems the Regents were trying to avoid; printing the extreme thins also required the use of unusually high-quality paper. The 38 glyphs of the 1868–69 book font were basically the same as the 38 glyphs used in printing articles in the *Deseret News* in 1859–60 and 1864; the only significant difference was that the old 3 glyph was mirror-imaged to ε.

³² Journal History, 16 February 1859.

peep	ᵛ	ᵛ	bib
tot	1	ᵛ	dead
kick	ᵛ	ᵛ	gag
fee	J	ᵛ	vow
thigh	ᵛ	ᵛ	they
so	ᵛ	ᵛ	zoo
sure	ᵛ	ᵛ	measure
church	ᵛ	ᵛ	judge
yea	\	/	*woe
hung	ᵛ	ᵛ	ha-ha
loll	c	ᵛ	roar
mime*	s	ᵛ	none
if	1	ᵛ	eat
egg	ᵛ	ᵛ	age
ash*	J	ᵛ	ice
ado*	ᵛ	ᵛ	up
on	ᵛ	0	oak
wool	v	ᵛ	ooze
out	<	>	oil
ah*	s	ᵛ	awe
are	ᵛ	ᵛ	or
air	ᵛ	ᵛ	urge
array	ᵛ	ᵛ	ear
ian	r	ᵛ	yew
the	ᵛ	ᵛ	of
and	\	1	to

*written top-down or right-left
 ▶ for proper names,
 use »Namer« dot (eg, ᵛor, Rome).

Figure 18: The Shaw or “Shavian” Alphabet was designed by typographer Kingsley Read and has inspired a number of other professional typographers, including Ross DeMeyere (<http://www.demeyere.com/shavian/>). The glyphs are simple and harmonious; ascenders and descenders give words distinctive shapes and avoid monotony. Copyright © 2002 DeMeyere Design Incorporated. All rights reserved. Reproduced by permission.

אורכ שן גרזלזלזקד

1. טגתא אדכ, קטא ק שן 1 אורכ זק. אן.
 אססמ ט זלזלז א ק זלז. אטקא אר
 אגד אסון אסל אורכ ז אן קטא, אטלזא
 אטלזא ז ארז קדלזק ז אטלז א.
 אט אט אן א טזל אורכ ז א קא.

Figure 19: In this extract of Shavian script, the title is set in the Ghoti (pronounced “fish”) font, and the body in the Androcles font, both by Ross DeMeyere. Copyright © 2002 DeMeyere Design Incorporated. All rights reserved. Reproduced by permission.

of their uniformity. Another objection some have urged against them has been that they are entirely new, and we should have characters as far as possible with which we are familiar: and they have felt that we should use them as far as they go and adopt new characters only for the sounds which our present letters do not represent.

There is a system known as the [Benn] Pitman system of phonetics which possesses the advantages alluded to. Mr. Pitman has used all the letters of the alphabet as far as possible and has added seventeen new characters to them, making an alphabet of forty-three letters. The Bible, a dictionary and a number of other works, school books, etc., have been printed in these new characters, and it is found that a person familiar with our present method of reading can learn in a few minutes to read those works printed after this system. We think it altogether likely that the regents of the University will upon further examination adopt this system for use in this Territory.

So while the Deseret Alphabet was dead, the Mormons hadn’t yet given up on spelling reform. In July of 1877, Orson Pratt was sent to Liverpool to arrange to have *The Book of Mormon* and *The Book of Doctrine and Covenants*, another book of Mormon scripture, printed in the Benn Pitman orthography, “with the exception of two or three characters”.³³

³³ *Journal of Discourses*, vol. XIX, p. 112.

But in August of that year, after most of the specially ordered phonotype had arrived from London, Brigham Young died; Orson Pratt was called back home, and the Mormons never dabbled in orthographical reform again.

It has been written, and repeated numerous times, that “the Deseret Alphabet died with Brigham Young”; however, the Deseret Alphabet had already been dead for at least a couple of years, and what died with Brigham Young was a very serious project, well in progress, to print Mormon scripture in a slight modification of Benn Pitman’s “American phonotypy”.

4 The Deseret Alphabet in Unicode

4.1 The Character Inventory and Glyphs

The Deseret Alphabet was first added to the Unicode 3.1 standard³⁴ in 2001, in the surrogate space 10400–1044F, mostly through the efforts of John H. Jenkins of Apple Computer.³⁵ It holds some distinction as the first script proposed for the surrogate space; as Jenkins describes it, “Nobody started to implement surrogates because there were no characters using them, and nobody wanted their characters to be encoded using surrogates because nobody was implementing them.”³⁶ The Deseret Alphabet, being a real but pretty dead script, was chosen as a pioneer—or sacrificial lamb—to break the vicious circle.

The Unicode 3.1 encoding handled only the 38-letter version of the Deseret Alphabet (this made 76 characters, including uppercase and lowercase) used in the printed books of 1868–69. The implementors were honestly unaware that earlier 39- and 40-letter versions of the Alphabet had been seriously used, and so might need to be encoded. I later argued vigorously³⁷ for the addition of the /^j/ and /^u/ letters used in several earlier versions of the Alphabet, including the one used in the Haskell journal and Shelton letters that I have transcribed. John Jenkins backed me up³⁸ and again deserves the credit for dealing with most of the paperwork and bureaucracy.

The two new letters were included in Unicode 4.0, but unfortunately I could not persuade them to use the 1859–60 glyphs \emptyset and ? as the citation glyphs; instead they went all the way back to the primitive glyphs of the 1854–55 charts. Unicode

	1040	1041	1042	1043	1044
0	\emptyset 10400	? 10410	S 10420	L 10430	D 10440
1	E 10401	7 10411	? 10421	L 10431	P 10441
2	E 10402	B 10412	L 10422	7 10432	E 10442
3	O 10403	7 10413	O 10423	9 10433	L 10443
4	O 10404	E 10414	7 10424	? 10434	? 10444
5	O 10405	C 10415	H 10425	E 10435	S 10445
6	7 10406	9 10416	? 10426	W 10436	E 10446
7	L 10407	D 10417	O 10427	? 10437	D 10447
8	L 10408	D 10418	\emptyset 10428	? 10438	S 10448
9	L 10409	P 10419	E 10429	7 10439	? 10449
A	7 1040A	E 1041A	E 1042A	B 1043A	L 1044A
B	9 1040B	L 1041B	O 1042B	7 1043B	O 1044B
C	? 1040C	? 1041C	O 1042C	E 1043C	7 1044C
D	E 1040D	S 1041D	O 1042D	C 1043D	H 1044D
E	W 1040E	E 1041E	7 1042E	9 1043E	? 1044E
F	? 1040F	D 1041F	L 1042F	D 1043F	O 1044F

Figure 20: The Deseret Alphabet as it appears in Unicode 4.0. Copyright © 1991–2003 Unicode, Inc. All rights reserved. Reproduced by permission of Unicode, Inc.

³⁴ <http://www.unicode.org/>

³⁵ <http://homepage.mac.com/jenkins/>

³⁶ <http://homepage.mac.com/jenkins/Deseret/Unicode.html>;
<http://homepage.mac.com/jenkins/Deseret/Computers.html>

³⁷ Unicode discussion document N2474 2002-05-17.

³⁸ Unicode discussion document N2473 2002-05-17.

Table 1: The Deseret Alphabet was added to Unicode by General American English speakers who honestly misunderstood the ɹ (/ɒ/) and 0 (/ɔ/) vowels, which have collapsed to /ɑ/ in their dialect, and renamed them confusingly as SHORT AH and LONG AH.

Char.	IPA	Original Name	Unicode Name
∂	/i/	e as in eat	LONG I
ε	/e/	a as in ate	LONG E
∂	/ɑ/	ah as in art	LONG A
0	/ɔ/	aw as in aught	LONG AH
0	/o/	o as in oat	LONG O
∅	/u/	oo as in ooze	LONG OO
†	/ɪ/	i as in it	SHORT I
ɹ	/ɛ/	e as in et	SHORT E
ɹ	/æ/	a as in at	SHORT A
ɹ	/ɒ/	o as in ot	SHORT AH
ɹ	/ʌ/	u as in ut	SHORT O
ɹ	/ʊ/	oo as in book	SHORT OO

fonts based on the current heterogeneous collection of glyphs will be useless for any practical typesetting of 40-letter Deseret Alphabet documents.

4.2 Unicode Character Names

The Unicode implementation of the Deseret Alphabet is also flawed by some changes to the letter names. Not to criticize anyone personally, but just for the record, there are several reasons why the name changes were ill-advised:

1. The Deseret Alphabet had a traditional set of letter names already established and available. Arbitrary changes in the names make it more difficult to compare the original charts and the Unicode charts.
2. Some early Deseret Alphabet writers, including George D. Watt, consciously or unconsciously confused the traditional letter names and their phonological values. Some of their spellings make sense only if the letters are read with their original names.
3. Some letter-name changes were made because the implementors simply did not hear and understand some of the vowel distinctions provided in the Deseret Alphabet; they were speakers of General American English, a dialect that has lost some of the vowel distinctions still present in English and New England dialects.

The last point is the most unfortunate. Consider Table 1: The original name for the Deseret 0 letter, which is /ɑ/ in IPA, was “ah”, using a com-

mon convention in English romanization whereby “ah” represents an unrounded low-back vowel. Most English speakers use this vowel in the words *father*, *bah* and *hah*. In England, and in much of New England, this vowel is distinct from the first vowel in *bother*, represented in Deseret Alphabet as ɹ or in IPA as /ɒ/, which is a rounded low-back vowel; thus for these speakers the words *father* and *bother* do not rhyme. But the rounded /ɒ/ has collapsed into unrounded /ɑ/ in General American English, so the words do rhyme for most Americans. Similarly, the Deseret 0 letter, IPA /ɔ/, represents a mid-low back rounded vowel that has also collapsed into /ɑ/ for many American speakers. It can still be heard quite distinctly in the speech of many New Yorkers, Philadelphians, and New Englanders in general. The original Deseret name for the 0, “aw”, used a common convention for representing this rounded vowel, which occurs in words like *law*, *flaw*, *paw*, *aught*, *caught*, etc. The equivalent letter in the Shaw Alphabet is appropriately named AWE. Not understanding the phonological distinctions involved, the implementors of Unicode renamed ɹ as SHORT AH and 0 as LONG AH, giving precisely the wrong clues to the pronunciation of these rounded vowels. Unfortunately, Unicode policy values consistency over accuracy, and it’s almost impossible to change character names once they have been adopted.

5 Digital Fonts for the Deseret Alphabet

5.1 Non-METAFONT Fonts

5.1.1 Kearney’s Deseret Font

A number of digital fonts have been designed for the Deseret Alphabet, most of them based on the 38-letter inventory and glyphs of the book font of 1868–69. The following is a very preliminary survey of fonts that I was able to find and test in early 2004.³⁹

The prize for the first digital font would seem to go to Greg Kearney, whose Deseret font was created about 1991 using Fontographer. Kearney (personal communication) says that his font, now in the public domain, was created for the LDS Church History Department, now the LDS Church Archives, as a display font for an exhibit.

I had difficulty testing this font⁴⁰ to input specific texts on my Mac OS X system, but see Figure 21

³⁹ The world of fonts, and especially amateur fonts, is woefully lacking in documentation. I would be extremely grateful for corrections and additions to the information in this section.

⁴⁰ <http://www.fontage.com/pages/deseret.html>; <http://funsite24.com/fo/d/>

as I can tell, it is no longer available, and numerous Internet links to Zion pages are dead. A webpage⁴⁶ dedicated to Deseret Alphabet fonts is a virtual cemetery of dead links.

5.1.5 Kass’s Code2001 Font

The freely available Code2001 font⁴⁷ by James Kass is a Plane 1 Unicode-based font, providing glyphs for the characters in the surrogate space, including Old Persian Cuneiform, Deseret, Tengwar, Cirith, Old Italic, Gothic, etc. Kass informs me that the glyphs (see Figure 25) were designed from scratch and resided originally in the Private Use Area of the Code2000 font until Deseret was officially accepted and assigned code points in the surrogate space.

5.1.6 Thibeault’s Deseret and Bartok’s HuneyBee Fonts

Daniel Thibeault took the Deseret Alphabet glyphs from the Code2001 font and transposed them into the ANSI range to make yet another font named Deseret.⁴⁸ Stephen Bartok’s HuneyBee font⁴⁹ was created in September 2003 by rearranging the glyphs in Thibeault’s Deseret font to effect a different keyboard layout (personal communication). In both fonts the glyphs are ultimately from the Code2001 font, already illustrated in Figure 25.

5.1.7 Elzinga’s Brigham Font

Dirk Elzinga of the Department of Linguistics and English Language at Brigham Young University is working on a new font called Brigham (see Figure 26), using FontForge, that is largely mono-width but judiciously uses thinner strokes for the loops.

5.1.8 Robertson’s Fonts

Graphic designer Christian Robertson is working on two fonts, “trying to make the Deseret Alphabet look good in type” (personal communication), which is quite a challenge. In his first font, Robertson is not afraid to “take out some of the curly queues that really mucked things up”, to rethink the representation of the short vowels, to add serifs, and even to introduce something like ascenders. The sample in Figure 27, kindly provided by Robertson, does not represent the latest version of his font, and the text is gibberish, but it illustrates his innovative approach. Robertson’s next font will be even more

challenging, designed for typesetting the early curative manuscripts from 1854–55.

5.2 Beesley’s METAFONT desalph Font and L^AT_EX Package

My own desalph font (see Figure 28) was created with METAFONT for the specific purpose of typesetting 40-letter Deseret Alphabet manuscripts from 1859–60. These documents were typically written with narrow nib pens, producing some thick-thin distinction, so the coding relies heavily on METAFONT penstroke commands. I took my inspiration from the pre-book charts of 1854–55, and from real handwriting. The penstrokes follow the path used to draw the glyphs, giving a hint of the original handwriting that is completely obscured in the Bodonified book font of 1868–69.

The desalph font is made available in a desalph package, which can be used in a L^AT_EX document much like the TIPA package.⁵⁰ The input of Deseret Alphabet characters can be done somewhat clumsily using commands like `\dalclongi` (Deseret Alphabet lowercase long i) for \mathfrak{a} or `\dauclongi` (Deseret Alphabet uppercase long i) for \mathfrak{A} . Inside `\textda{}` commands, a more convenient system of transliteration “shortcuts” can be used. As I was already somewhat comfortable with the shortcuts of the TIPA package, for entering IPA letters, I laid out the desalph font internally so that the same shortcuts could be used wherever possible. Simple commands were defined to enter diphthongs and affricates, which have no shortcuts in TIPA. A simply defined `\ipa{}` command allows the same commands to be used to enter equivalent IPA diphthongs and affricates. The principal entry commands are summarized in Table 2, and some extra commands for unusual and idiosyncratic glyphs are shown in Table 3. Uppercase letters, found in Deseret Alphabet but not in IPA, can be entered with corresponding uppercase “uc” commands with names like `\dauclongi`, or by placing the shortcut in the `\uc{}` command, e.g. `\uc{i}`.

The use of METAFONT allowed me to define the proper glyphs for the 1859–60 manuscripts, especially the \mathfrak{O} used for $/\mathfrak{o}^j/$ and the $\mathfrak{?}$ used for $/\mathfrak{i}u/$, which I have never seen in a printed chart or document.⁵¹ When I found a manuscript with the experimental new letter \mathfrak{l} for the neutral vowel called schwa ($/\mathfrak{a}/$), making a 41-letter alphabet, adding it to my METAFONT font was a simple exercise.

⁴⁶ <http://cgm.cs.mcgill.ca/~luc/deseret.html>

⁴⁷ <http://home.att.net/~jameskass/code2001.htm>

⁴⁸ <http://www.angelfire.com/pq/Urhixidur/Fonts/Fonts.html>

⁴⁹ <http://home.earthlink.net/~slbartok/projects/fonts.htm>

⁵⁰ <http://tooyoo.l.u-tokyo.ac.jp/~fkr/>

⁵¹ An \mathfrak{O} punch appears in the set of St. Louis punches of 1857, but it was not used when printing finally started in 1859.

ᐃᐅᐅᐅᐅᐅ ᐺᐺᐺᐺᐺ ᐃᐅ ᐃᐺᐺ ᐺᐺᐺᐺᐺᐺᐺᐺᐺᐺᐺᐺᐺᐺᐺᐺ
 ᐃᐅᐅᐅᐅᐅ ᐺᐺᐺᐺᐺ ᐃᐅ ᐃᐺᐺ ᐺᐺᐺᐺᐺᐺᐺᐺᐺᐺᐺᐺᐺᐺᐺᐺ

ᐺᐺᐺᐺᐺᐺᐺ, ᐺᐺ ᐺᐺ ᐺᐺᐺᐺᐺᐺᐺᐺ ᐺᐺ ᐺ ᐺᐺᐺᐺᐺ ᐺᐺ ᐺ ᐺᐺᐺ ᐺᐺ ᐺᐺᐺᐺ, ᐺᐺ ᐺᐺ ᐺᐺ ᐺᐺ ᐺᐺᐺᐺᐺᐺᐺ; ᐺᐺᐺ ᐺᐺ ᐺ ᐺᐺᐺᐺᐺᐺᐺ, ᐺᐺ ᐺᐺ ᐺ ᐺᐺᐺᐺᐺᐺ ᐺᐺ ᐺ ᐺᐺᐺ ᐺᐺ ᐺᐺᐺᐺᐺᐺ; ᐺᐺ ᐺᐺᐺ ᐺᐺ ᐺᐺᐺ ᐺᐺᐺᐺᐺᐺ: ᐺᐺᐺ ᐺᐺ ᐺᐺ ᐺᐺ ᐺᐺᐺᐺᐺᐺᐺᐺᐺ, ᐺᐺ ᐺᐺᐺ ᐺᐺ ᐺ ᐺᐺᐺᐺᐺᐺ ᐺᐺ ᐺᐺᐺᐺᐺᐺ ᐺᐺ ᐺᐺ ᐺᐺᐺᐺᐺᐺᐺᐺᐺᐺ. ᐺᐺᐺ ᐺᐺ ᐺᐺᐺ ᐺᐺ, ᐺᐺ ᐺᐺ ᐺᐺ ᐺᐺᐺ ᐺ ᐺᐺᐺ, ᐺᐺᐺ ᐺᐺ ᐺᐺᐺ ᐺᐺ ᐺ ᐺᐺᐺᐺᐺᐺᐺ.

Figure 26: Elzinga’s Brigham font.

ᐺᐺᐺᐺ ᐺᐺ ᐺᐺᐺᐺ ᐺ ᐺᐺ ᐺᐺ ᐺᐺᐺ ᐺᐺᐺᐺ ᐺᐺᐺ ᐺᐺ ᐺᐺᐺ
 ᐺᐺᐺᐺᐺ ᐺᐺᐺᐺᐺᐺ ᐺᐺ ᐺᐺᐺ ᐺᐺᐺ ᐺᐺᐺᐺᐺ ᐺᐺ ᐺᐺᐺ
 ᐺᐺᐺᐺᐺ ᐺᐺ ᐺᐺᐺ ᐺᐺᐺᐺ ᐺᐺ ᐺᐺ ᐺᐺᐺᐺ ᐺᐺ ᐺᐺᐺ
 ᐺᐺ ᐺᐺᐺᐺᐺᐺ ᐺᐺᐺ ᐺᐺ ᐺᐺᐺ ᐺᐺᐺᐺ ᐺᐺᐺᐺ ᐺᐺ ᐺᐺᐺ
 ᐺᐺ ᐺᐺᐺᐺ ᐺᐺ ᐺᐺᐺ ᐺᐺ ᐺᐺᐺ ᐺᐺᐺᐺ ᐺᐺ ᐺᐺᐺᐺ ᐺᐺ
 ᐺᐺᐺᐺ ᐺᐺ ᐺᐺᐺᐺᐺᐺ ᐺᐺᐺ ᐺᐺ ᐺᐺᐺ ᐺᐺᐺᐺ ᐺᐺᐺᐺᐺ
 ᐺᐺ ᐺᐺᐺ ᐺᐺᐺᐺ ᐺᐺ ᐺᐺ ᐺᐺᐺᐺᐺ ᐺᐺ ᐺᐺᐺᐺ ᐺᐺ ᐺᐺᐺᐺ
 ᐺᐺᐺᐺ ᐺᐺ ᐺᐺᐺᐺ ᐺᐺ ᐺᐺᐺᐺ ᐺᐺ ᐺᐺᐺᐺ ᐺᐺᐺᐺ ᐺᐺᐺᐺᐺ

Figure 27: Robertson’s experimental font.

ᐃᐅᐅᐅᐅᐅ ᐺᐺᐺᐺᐺ ᐃᐅ ᐃᐺᐺ ᐺᐺᐺᐺᐺᐺᐺᐺᐺᐺᐺᐺᐺᐺᐺᐺ ᐅᐺ
 ᐃᐅᐅᐅᐅᐅ ᐺᐺᐺᐺᐺ ᐃᐅ ᐃᐺᐺ ᐺᐺᐺᐺᐺᐺᐺᐺᐺᐺᐺᐺᐺᐺᐺᐺ ᐅᐺ

ᐺᐺᐺᐺᐺᐺᐺ, ᐺᐺ ᐺᐺ ᐺᐺᐺᐺᐺᐺᐺᐺ ᐺᐺ ᐺ ᐺᐺᐺᐺᐺ ᐺᐺ ᐺ ᐺᐺᐺ ᐺᐺ ᐺᐺᐺᐺ, ᐺᐺ ᐺᐺᐺ ᐺᐺ ᐺ ᐺᐺᐺᐺᐺᐺᐺ; ᐺᐺᐺ ᐺᐺ ᐺ ᐺᐺᐺᐺᐺᐺᐺ, ᐺᐺ ᐺᐺ ᐺ ᐺᐺᐺᐺᐺᐺ ᐺᐺ ᐺ ᐺᐺᐺ ᐺᐺ ᐺᐺᐺᐺᐺᐺ; ᐺᐺ ᐺᐺᐺ ᐺᐺ ᐺᐺᐺ ᐺᐺᐺᐺᐺᐺᐺ: ᐺᐺᐺ ᐺᐺ ᐺᐺ ᐺᐺ ᐺᐺᐺᐺᐺᐺᐺᐺᐺ, ᐺᐺ ᐺᐺᐺ ᐺᐺ ᐺ ᐺᐺᐺᐺᐺᐺ ᐺᐺ ᐺᐺᐺᐺᐺᐺᐺ ᐺᐺ ᐺᐺ ᐺᐺᐺᐺᐺᐺᐺᐺᐺᐺᐺ. ᐺᐺᐺ ᐺᐺ ᐺᐺᐺ ᐺᐺ, ᐺᐺ ᐺᐺ ᐺᐺ ᐺᐺᐺ ᐺ ᐺᐺᐺᐺ, ᐺᐺ ᐺᐺ ᐺᐺᐺ ᐺᐺ ᐺ ᐺᐺᐺᐺᐺᐺᐺᐺᐺ.

Figure 28: Beesley’s desalph font.

Table 3: Extra commands used to enter rare and idiosyncratic Deseret Alphabet glyphs.

ᐅ	\daucslju	St. Louis 1857 font, unused glyph for / ^l u/
ᐅ	\dauchaskoi	Haskell’s idiosyncratic glyph for / ^o j/
	\daucschwa	Shelton’s proposed glyph for schwa /ə/
ᐺ	\daucspellerow	Deseret Phonetic Speller glyph for /a ^w /


```

\documentclass[] {article}
\usepackage{times}
\usepackage{desalph}
\usepackage{tipa}
% commands used in \ipa{}, parallel to commands in \textda{}, to get
% an equivalent phonemic IPA transliteration of Deseret Alphabet
\newcommand{\ipa}[1]{\{\tipaencoding%
\providecommand{\aI}{\renewcommand{\aI}{a\textsuperscript{j}\xspace}%
\providecommand{\aJ}{\renewcommand{\aJ}{a\textsuperscript{j}\xspace}%
\providecommand{\OI}{\renewcommand{\OI}{O\textsuperscript{j}\xspace}%
\providecommand{\OJ}{\renewcommand{\OJ}{O\textsuperscript{j}\xspace}%
\providecommand{\aU}{\renewcommand{\aU}{a\textsuperscript{w}\xspace}%
\providecommand{\aW}{\renewcommand{\aW}{a\textsuperscript{w}\xspace}%
\providecommand{\ju}{\renewcommand{\ju}{\textsuperscript{j}u\xspace}%
\providecommand{\Ju}{\renewcommand{\Ju}{\textsuperscript{j}u\xspace}%
\providecommand{\dZ}{\renewcommand{\dZ}{\textdyoghlig\xspace}%
\providecommand{\tS}{\renewcommand{\tS}{\texttshlig\xspace}#1}}

\begin{document}
\begin{center}
A sample of Deseret Alphabet entered using shortcuts:\\
\textda{i e A O o u I E \ae{} 6 2 U \aI{} \OI{} \aU{} \ju{}
w j h p b t d \tS{} \dZ{} k g f v T D s z S Z r l m n N}
\smallskip
Parallel phonemic IPA entered using the same shortcuts:\\
\ipa{i e A O o u I E \ae{} 6 2 U \aI{} \OI{} \aU{} \ju{}
w j h p b t d \tS{} \dZ{} k g f v T D s z S Z r l m n N}
\end{center}
\end{document}

```

Figure 29: A skeleton L^AT_EX example using the TIPA and desalph packages

I discovered an uncatalogued and unidentified “Indian Vocabulary” in the LDS Church Archives, and I was able to identify it as English-to-Hopi. I have argued [8] that it was written by Marion J. Shelton during this mission, and it appears to be the oldest written evidence of the Hopi language.

The entire vocabulary has now been typed into an XML format, with fields added for modern English and Hopi orthography, modern dictionary definitions, and comments and references of various kinds. The XML file is downtranslated using a Perl-language script, with the helpful Perl XML::Twig package,⁵⁵ to produce L^AT_EX source code with Deseret Alphabet output, using the `desalph` package and font, and equivalent phonemic IPA output, using the TIPA package. The use of XML, the `desalph` font, TIPA and L^AT_EX allows me and my co-author Dirk Elzinga to reproduce this extraordinary document for study and publication. Creating and maintaining the original data in an XML format gives us all the advantages of XML validation and abstraction; and the flexibility of downtranslation to L^AT_EX

allows us to format the output in different ways suitable for proofreading or for final publication.

The English-Hopi Vocabulary (see Figure 30) is written entirely in the Deseret Alphabet and includes 486 entries like the following

†‡‡††‡††‡‡   ‡‡C.‡‡.†‡

with an English word on the left and a Hopi word in Third Mesa (Orayvi) dialect on the right. Encoded as XML, and with auxiliary information added, this entry appears as shown in Figure 31. The XML file is validated using a Relax NG schema. Downtranslation of the XML entry currently yields the L^AT_EX output in Figure 32, which is a line in a table. When typeset, the entry appears as shown in Table 4. This open tabular format is ideal for proofreading, and for the final paper all that will be required is a modified Perl script to downtranslate the same XML file into other L^AT_EX codes that waste less space.

I have also transcribed the journal of Thales H. Haskell, kept in the Deseret Alphabet from October through December of 1859, and will include it in a general history of the second mission to the Hopi [7]. Here, for reading practice, is an extract from

⁵⁵ <http://www.xmltwig.com/xmltwig/>


```

340 & \raggedright \index{rabbit stick, 340} rabbit stick \\
\textda{r\ae{}bIt-stIk} \\
\ipa{r\ae{}bIt-stIk} & \raggedright \ipa{pe\tS{}}.ko.ho} \\
\textda{pe\tS{}}.ko.ho} & HD p.\@ 449: puts$|koho
‘rabbit stick, a flat boomerang-like stick
used for hunting; used for throwing and hitting
it on the run’\

```

Figure 32: L^AT_EX output from downloading an XML entry.

Table 4: Entry of the English-Hopi vocabulary typeset for proofreading.

340	rabbit stick †ŋɪːsɪŋ ræbit-stik	peŋ.ko.ho ʔɪC.ŋŋ.ʔŋ	HD p. 449: puts koho ‘rabbit stick, a flat boomerang-like stick used for hunting; used for throwing and hitting it on the run’
-----	---------------------------------------	------------------------	--

\$7000 or \$8000, there has always been some interest in retypesetting it.

The Deseret First Book and *The Deseret Second Book* had print runs of 10,000 copies each, are therefore much more plentiful, and copies today go for around \$200. *The Deseret First Book* has even been reprinted photographically for sale to tourists as a Utah curiosity [26], and the text has been keyed in by John Jenkins, and proofread by Michael Everson and by myself. Such projects are of interest to linguists who want to search the texts electronically.

In 1967, LDS Church archivists found a bundle of forgotten Deseret Alphabet manuscripts, some of them ready for the typesetter but never printed [32]. These include *The Doctrine and Covenants*, with the *Lectures on Faith*; the *Catechism* of John Jaques; and the entire text of the *Bible*. The LDS Church Archives also hold the *History of Brigham Young*, a number of letters, an unfinished *Deseret Phonetic Speller*, journals, letters and probably a number of other documents still to be found.

7 Conclusion

Although the Deseret Alphabet was never intended for secrecy [6], few people then or now can be persuaded to learn it, and a number of interesting documents have been ignored and unstudied for over 140 years. The letters and journals are of interest to historians, and the phonemically written texts are also of interest to linguists. With the help of XML, L^AT_EX, TIPA and new digital fonts for the Deseret Alphabet, these neglected documents are coming to light again.

References

- [1] Douglas D. Alder, Paula J. Goodfellow, and Ronald G. Watt. Creating a new alphabet for Zion: The origin of the Deseret Alphabet. *Utah Historical Quarterly*, pages 275–286, 1984.
- [2] Alfred Baker. *The Life of Sir Isaac Pitman: Inventor of Phonography*. Isaac Pitman and Sons, New York, 1908.
- [3] Edward Bateman. A brief history of the Deseret Alphabet. *Emigre*, (52):72–77, 1999. Fall.
- [4] Kenneth R. Beesley. The Deseret Alphabet: Can orthographical reform for English succeed? 1975.
- [5] Kenneth R. Beesley. Dialect determinants in the Deseret Alphabet journal of Thales H. Haskell. *Deseret Language and Linguistic Society Bulletin*, (3):2–35, 1977.
- [6] Kenneth R. Beesley. The Deseret Alphabet in Unicode. In *Proceedings of the 22nd International Unicode Conference*, volume 2, San Jose, California, September 11–13 2002. Unicode Consortium. Paper C10.
- [7] Kenneth R. Beesley. The second Mormon mission to the Hopi: 1859–60. Forthcoming, 2004.
- [8] Kenneth R. Beesley and Dirk Elzinga. An 1860 English-to-Hopi vocabulary written in the Deseret Alphabet. Forthcoming, 2004.
- [9] Juanita Brooks. Journal of Thales H. Haskell. *Utah Historical Quarterly*, pages 69–98, 1944.
- [10] Bernard Desgraupes. *METAFONT: Guide pratique*. Vuibert, Paris, 1999.
- [11] Leah R. Frisby and Hector Lee. The Deseret readers. *Utah Humanities Review*, 1:240–244, 1947.

- [12] IPA. *Handbook of the International Phonetic Association: A Guide to the Use of the International Phonetic Alphabet*. Cambridge University Press, Cambridge, 1999.
- [13] Harry C. James. *The Hopi Indians: Their history and their culture*. Caxton, Caldwell, ID, 1956.
- [14] Harry C. James. *Pages from Hopi History*. The University of Arizona Press, Tucson, AZ, 1974.
- [15] J. Kelly. The 1847 Alphabet: An episode of phonotypy. In R. E. Asher and Eugénie J. A. Henderson, editors, *Towards a History of Phonetics*, pages 248–264. Edinburgh University Press, Edinburgh, 1981.
- [16] Donald E. Knuth. *The METAFONTbook*. Addison-Wesley, New York, 1986.
- [17] Peter Ladefoged. *A Course in Phonetics*. Harcourt College Publishers, Orlando, FL, fourth edition, 2001.
- [18] P. A. D. MacCarthy. The Bernard Shaw alphabet. In Werner Haas, editor, *Alphabets for English*, pages 105–117. Manchester University Press, Manchester, 1969.
- [19] Samuel C. Monson. The Deseret Alphabet. Master’s thesis, Columbia University, 1947.
- [20] Samuel C. Monson. The Deseret Alphabet. In *Utah Academy of Sciences, Arts & Letters*, volume 30, pages 1952–53, 1953.
- [21] William V. Nash. The Deseret Alphabet. Master’s thesis, University of Illinois Library School, Urbana, Illinois, May 1957.
- [22] Douglas Allen New. *History of the Deseret Alphabet and other attempts to reform English Orthography*. PhD thesis, Utah State University, 1985.
- [23] Charles S. Peterson. The Hopis and the Mormons: 1858-1873. *Utah Historical Quarterly*, 39(2):179–194, Spring 1971.
- [24] Benn Pitman. *Sir Isaac Pitman: His Life and Labors*. C. J. Krehbiel, Cincinnati, 1902.
- [25] Thomas Allen Reed. *A Biography of Isaac Pitman: Inventor of Phonography*. Griffith, Farran, Okeden and Welsh, London, 1890.
- [26] Regents of the Deseret University. *Deseret Alphabet: The Deseret First Book*. Buffalo River Press, Salt Lake City, historical reprint edition, 1996.
- [27] Jules Rémy and Julius Brenchley. *A Journey to Great-Salt-Lake City*. W. Jeffs, London, 1861.
- [28] Bernard Shaw. *Androcles and the Lion*. Penguin, Harmondsworth, Middlesex, Shaw Alphabet edition, 1962.
- [29] Albert E. Smith. Thales Hastings Haskell: Pioneer, scout, explorer, Indian missionary. Typescript, Brigham Young University Library, 1964.
- [30] Hosea Stout. *Journal of Hosea Stout*. University of Utah Press, Salt Lake City, 1964.
- [31] Roby Wentz. *Thirty-Eight Mormon Characters: A Forgotten Chapter in Western Typographic History: For the Zamorano Club Jubilee*. Zamorano Club, Los Angeles, 1978.
- [32] Albert L. Zobell. Deseret Alphabet manuscripts found. *Improvement Era*, pages 10–11, July 1967.

Calendar

2004

- | | |
|---|---|
| <p>Mar 22–
May 7 In Flight: A traveling juried exhibition of books by members of the Guild of Book Workers. University of Washington, Seattle, Washington. Sites and dates are listed at http://palimpsest.stanford.edu/byorg/gbw.</p> <p>Apr 5–8 Book History Workshop, Institute d’histoire du livre, Lyon, France. For information, visit http://ihl.enssib.fr/.</p> <p>Apr 19–21 Seybold Seminars Amsterdam 2004, Netherlands. For information, visit http://www.seybold365.com/ams2004/.</p> <p>Apr 29–
May 1 Wells Book Arts Center symposium, Matter & Spirit: The Genesis & Evolution of the Book, Aurora, New York. For information, visit http://aurora.wells.edu/~wbac/bookarts/events.html.</p> <p>Apr 30–
May 3 BachoT_EX 2004, 12th annual meeting of the Polish T_EX Users’ Group (GUST), Bachotek, Brodnica Lake District, Poland. For information, visit http://www.gust.org.pl/BachoTeX/2004/.</p> <p>May 20–
Jul 7 In Flight: A traveling juried exhibition of books by members of the Guild of Book Workers. Emory University, Atlanta, Georgia. Sites and dates are listed at http://palimpsest.stanford.edu/byorg/gbw.</p> <p>Jun 11–16 ALLC/ACH-2004, Joint International Conference of the Association for Computers and the Humanities, and Association for Literary and Linguistic Computing, “Computing and Multilingual, Multicultural Heritage”, Göteborg University, Sweden. For information, visit http://www.hum.gu.se/allcach2004/ or the organization web site at http://www.ach.org.</p> | <p>Jun 14 GUTenberg Journée L^AT_EX, Paris, France. For information, visit http://www.gutenberg.eu.org/manifestations/gut2004/.</p> <p>Jul 5–
Aug 6 Rare Book School, University of Virginia, Charlottesville, Virginia. Many one-week courses on topics concerning typography, bookbinding, calligraphy, printing, electronic texts, and more. For information, visit http://www.virginia.edu/oldbooks.</p> <p>Jun 24–29 2nd International Conference on Typography and Visual Communication: Communication and new technologies, Thessaloniki, Greece. For information, visit http://www.uom.gr/uompress/.</p> <p>July 16–
Aug 28 In Flight: A traveling juried exhibition of books by members of the Guild of Book Workers. Columbia College, Chicago, Illinois. Sites and dates are listed at http://palimpsest.stanford.edu/byorg/gbw.</p> <p>Jul 19–22 Practical T_EX 2004, San Francisco, California. A user-oriented conference sponsored by TUG. For information, visit http://www.tug.org/practicaltex2004/.</p> <p>Jul 20–24 SHARP Conference (Society for the History of Authorship, Reading and Publishing), Lyon, France For information, visit http://sharpweb.org/.</p> <p>Jul 22–25 TypeCon2004, “Type High”, San Francisco, California. For information, visit http://www.typecon2004.com/.</p> <p>Aug 2–6 <i>Extreme</i> Markup Languages 2004, Montréal, Québec. For information, visit http://www.extrememarkup.com/extreme/.</p> <p>Aug 14–17 International Conference on Computing, Communications and Control Technologies, University of Texas, Austin, Texas. For information, visit http://www.iiisci.org/ccct2004/website/.</p> |
|---|---|

Status as of 1 April 2004

For additional information on TUG-sponsored events listed here, contact the TUG office (+1 503 223-9994, fax: +1 503 223-3960, e-mail: office@tug.org). For events sponsored by other organizations, please use the contact address provided.

An updated version of this calendar is online at <http://www.tug.org/calendar/>.

Additional type-related events are listed in the Typophile calendar, at

<http://www.icalx.com/html/typophile/month.php?cal=Typophile>.

Aug 8–12 SIGGRAPH 2004, Los Angeles,
California. For information, visit
<http://www.siggraph.org/calendar/>.
Aug 16–20 Seybold San Francisco, San Francisco,
California. For information, visit <http://www.seybold365.com/sf2004/>.

Oct 9 First meeting of GuIT (Gruppo
utilizzatori Italiani di T_EX),
Pisa, Italy. For information, visit
[http://www.guit.sssup.it/
GuITmeeting/2004/2004.en.html](http://www.guit.sssup.it/GuITmeeting/2004/2004.en.html).
Oct 18–19 Third Annual St. Bride Conference,
“Bad Type”, London, England.
For information, visit [http://
www.stbride.org/conference.html](http://www.stbride.org/conference.html).

TUG 2004

Democritus University of Thrace,
Xanthi, Greece.

Aug 30– The 25th annual meeting of the T_EX
Sep 3 Users Group, “XML and Digital
Typography”. For information, visit
<http://www.tug.org/tug2004/>.

Institutional Members

American Mathematical Society,
Providence, Rhode Island

Banca d'Italia,
Roma, Italy

Center for Computing Science,
Bowie, Maryland

Certicom Corporation,
Mississauga, Ontario, Canada

CNRS - IDRIS,
Orsay, France

CSTUG, *Praha, Czech Republic*

Duke University, Vesic Library,
Durham, North Carolina

Florida State University,
School of Computational Science
and Information Technology,
Tallahassee, Florida

IBM Corporation,
T J Watson Research Center,
Yorktown, New York

Institute for Advanced Study,
Princeton, New Jersey

Institute for Defense Analyses,
Center for Communications
Research, *Princeton, New Jersey*

KTH Royal Institute of
Technology, *Stockholm, Sweden*

Masaryk University,
Faculty of Informatics,
Brno, Czechoslovakia

Max Planck Institut
für Mathematik,
Bonn, Germany

National Association of
Mathematics Students,
Lagos, Nigeria

New York University,
Academic Computing Facility,
New York, New York

Princeton University,
Department of Mathematics,
Princeton, New Jersey

Siemens Corporate Research,
Princeton, New Jersey

Springer-Verlag Heidelberg,
Heidelberg, Germany

Stanford Linear Accelerator
Center (SLAC),
Stanford, California

Stanford University,
Computer Science Department,
Stanford, California

Stockholm University,
Department of Mathematics,
Stockholm, Sweden

University College, Cork,
Computer Centre,
Cork, Ireland

University of Delaware,
Computing and Network Services,
Newark, Delaware

Université Laval,
Ste-Foy, Québec, Canada

University of Oslo,
Institute of Informatics,
Blindern, Oslo, Norway

Uppsala University,
Uppsala, Sweden

Vanderbilt University,
Nashville, Tennessee



Promoting the use of
TeX throughout the
world.

mailing address:

P.O. Box 2311
Portland, OR 97208-2311 USA

shipping address:

1466 NW Naito Parkway
Suite 3141
Portland, OR 97209-2820 USA

phone: +1 503-223-9994
fax: +1 503-223-3960
email: office@tug.org
web: http://www.tug.org

President Karl Berry
Vice-President Kaja Christiansen
Treasurer Samuel Rhoads
Secretary Susan DeMeritt
Executive Director Robin Laakso

2004 TeX Users Group Membership Form

TUG membership rates are listed below. Please check the appropriate boxes and mail the completed form with payment (in US dollars, drawn on a US bank) to the mailing address at left. If paying by credit/debit card, you may alternatively fax the form to the number at left or join online at <http://tug.org/join.html>. The web page also provides more information than we have room for here.

Status (check one) New member Renewing member

	Rate	Amount
<input type="checkbox"/> Early bird membership for 2004 (TUGboat, software) <i>After May 31, dues are \$75.</i>	\$65	_____
<input type="checkbox"/> Special membership for 2004 (TUGboat, software) <i>You may join at this special rate (\$45 after May 31) if you are a senior (62+), student, new graduate, or from a country with a modest economy. Please circle the applicable choice. See http://tug.org/join.html for more information.</i>	\$35	_____
<input type="checkbox"/> Subscription for 2004 (TUGboat, software) non-voting	\$85	_____
<input type="checkbox"/> Institutional membership for 2004 (TUGboat, software) <i>Includes up to seven individual memberships.</i>	\$500	_____

Last year's materials

If you were not a TUG member in 2003, this is your option to receive software immediately.

<input type="checkbox"/> TeX Live 2003 software <i>2 CD's and 1 DVD which also includes CTAN.</i>	\$15	_____
<input type="checkbox"/> CTAN 2003 CD-ROMs	\$15	_____
<input type="checkbox"/> TUGboat Volume 24	\$15	_____

Voluntary donations

<input type="checkbox"/> General TUG contribution	_____
<input type="checkbox"/> Bursary Fund contribution <i>Financial assistance for attending the TUG Annual Meeting.</i>	_____
<input type="checkbox"/> TeX Development Fund contribution <i>Financial assistance for technical projects.</i>	_____
<input type="checkbox"/> Send me CTAN on CD (shipped on DVD to everyone)	n/a

Total \$ _____

Tax deduction: \$30 of the early bird membership fee is deductible, at least in the US.

Multi-year orders: To join for more than one year at this year's rate, just multiply.

Payment (check one) Payment enclosed Visa/MasterCard/AmEx

Account Number: _____

Exp. date: _____ Signature: _____

Privacy: TUG uses your personal information only to send products, publications, notices, and (for voting members) official ballots. TUG neither sells its membership list nor provides it to anyone outside of its own membership.

Electronic notices will generally reach you much earlier than printed ones. However, you may choose not to receive any email from TUG, if you prefer.

Do *not* send me any TUG notices via email.

Name _____

Department _____

Institution _____

Address _____

City _____ State/Province _____

Postal code _____ Country _____

Email address _____

Phone _____ Fax _____

Position _____ Affiliation _____

TEX Consulting & Production Services

Loew, Elizabeth

President, TEXniques, Inc.
675 Massachusetts Avenue, 6th Floor
Cambridge, MA 02139
(617) 876-2333; Fax: (781) 344-8158
Email: loew@texniques.com

Complete book and journal production in the areas of mathematics, physics, engineering, and biology. Services include copyediting, layout, art sizing, preparation of electronic figures; we keyboard from raw manuscript or tweak TEX files.

Ogawa, Arthur

40453 Cherokee Oaks Drive
Three Rivers, CA 93271-9743
(209) 561-4585
Email: arthur_ogawa@teleport.com

Bookbuilding services, including design, copyedit, art, and composition; color is my speciality. Custom TEX macros and L^AT_EX₂ ϵ document classes and packages. Instruction, support, and consultation for workgroups and authors. Application development in L^AT_EX, TEX, SGML, PostScript, Java, and C++. Database and corporate publishing. Extensive references.

Veytsman, Boris

2239 Double Eagle Ct.
Reston, VA 20191
(703) 860-0013
Email: borisv@lk.net

I provide training, consulting, software design and implementation for Unix, Perl, SQL, TEX, and L^AT_EX. I have authored several popular packages for L^AT_EX and `latex2html`. I have contributed to several web-based projects for generating and typesetting reports. For more information please visit my web page: <http://users.lk.net/~borisv>.

The Unicorn Collaborative, Inc., Ted Zajdel

115 Aspen Drive, Suite K
Pacheco, CA 94553
(925) 689-7442
Email: contact@unicorn-collab.com

We are a technical documentation company, initiated in 1990, which time, strives for error free, seamless documentation, delivered on time, and within budget. We provide high quality documentation services such as document design, graphic design and copy editing. We have extensive experience using tools such as FrameMaker, TEX, L^AT_EX, Word, Acrobat, and many graphics programs. One of our specialties is producing technical manuals and books using L^AT_EX and TEX. Our experienced staff can be trained to use any tool required to meet your needs. We can help you develop, rewrite, or simply copy-edit your documentation. Our broad experience with different industries allows us to handle many types of documentation including, but not limited to, software and hardware systems, communications, scientific instrumentation, engineering, physics, astronomy, chemistry, pharmaceuticals, biotechnology, semiconductor technology, manufacturing and control systems. For more information see our web page: <http://www.unicorn-collab.com>.

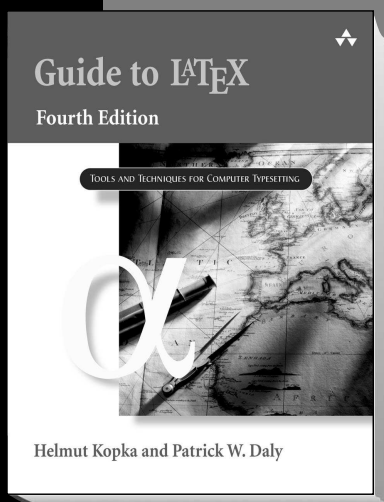
The information here comes from the consultants themselves. We do not include any information we know to be false, but we cannot check out any of the information; we are transmitting it to you as it was given to us and do not promise it is correct. Also, this is not an endorsement of the people listed here. We have no opinions and usually no information about the abilities of any specific person. We provide this list to enable you to contact service providers and decide for yourself whether to hire one.

The TUG office mentions the consultants listed here to people seeking TEX workers. If you'd like to be included, or run to a larger ad in *TUGboat*, please contact the office or see our web pages:

TEX Users Group
1466 NW Naito Parkway, Suite 3141
Portland, OR 97208-2311, U.S.A.
Phone: +1 503 223-9994
Fax: +1 503 223-3960
Email: office@tug.org
Web: <http://tug.org/consultants.html>
<http://tug.org/TUGboat/advertising.html>

Guide to L^AT_EX

Fourth Edition



HELMUT KOPKA AND PATRICK W. DALY

In this completely revised edition, the authors cover the L^AT_EX₂_ε standard and offer more details, examples, exercises, tips, and tricks. They go beyond the core installation to describe the key contributed packages that have become essential to L^AT_EX processing. *Guide to L^AT_EX, Fourth Edition*, will prove indispensable to anyone wishing to gain the benefits of L^AT_EX.

Available at fine bookstores every where.

ISBN: 0-321-17385-6

For more information, visit: www.awprofessional.com/titles/0321173856


Addison
Wesley

TUG 2005 Announcement and Call for Papers

TUG 2005 will be held in Wuhan, China during August 23–25, 2005. CTUG (Chinese T_EX User Group) has committed to undertake the conference affairs, and now announces the call for papers.

Why go to China for TUG 2005?

For fun!

This is the first TUG conference to be held in China. Wuhan is close to the birthplace of Taoism and the Three Gorges Reservoir. China is also the birthplace of typography in ancient times, and is simply a very interesting place to go.

For keeping up with the community!

The T_EX community in China has been growing over the years. China is one of the few countries in the world which has heavily applied free software (including T_EX, GNU/Linux, and more) in industry. The rich human resources and the creative T_EXackers have become a part of the engine driving the global T_EX community. TUG'05 is a good opportunity to meet them.

For your future!

The growing market is ready to use your expertise. Many libraries, publishing houses, and scientific organizations in China are eager to use your T_EX expertise.

Please submit abstracts for papers to tug2005@tug.org. For more information about TUG 2005, please visit:

<http://tug.org/tug2005>

Workshops and Presentations: \LaTeX , \TeX , \AMS-L\TeX , \ConTeXt

Practical \TeX 2004: Training and Techniques

Holiday Inn Fisherman's Wharf
San Francisco, California
July 19–22, 2004

<http://tug.org/practicaltex2004>
conferences@tug.org



Keynote address: *Peter Flynn, Silmaril Consultants*

Who should attend?

Mathematicians * University & corporate (L \TeX) \TeX documentation staff *
Students * Publishing company production staff * Scientists * Researchers

... and anyone who uses or is considering using the L \TeX and T \TeX technical documentation system.

Further information

This three-day conference will have user-oriented presentations and workshops on L \TeX , T \TeX , AMS-L \TeX , ConT \TeX t and their use in document production. Parallel tracks will be offered to accommodate a wide range of users, from beginning to advanced. Expert speakers and lecturers will present experiences and techniques useful to all.

Conference attendees will enjoy an opening night reception and a Chinatown banquet. Coffee and lunch will be served each day of the meeting. You will have easy access to many of San Francisco's colorful sights, just steps from the Fisherman's Wharf area.

Post-conference workshops

On the fourth day, July 22, courses will be offered focusing on specific areas: Intermediate and Advanced L \TeX training, Introduction to ConT \TeX t, and T \TeX on the Web.

Contact: conferences@tug.org

- Registration forms and hotel reservations are on the web site.
- If you'd like to promote your products and services, or simply support the conference, see the web site for sponsorship and advertising options.

Hope to see you there!

(Sponsored by the T \TeX Users Group.)

TeX Users Group

Preprints for the 2004 Annual Meeting:

TeX in the Era of Unicode

	2	Karl Berry / <i>Editorial remarks</i>
Electronic Documents	3	Christos KK Loverdos and Apostolos Syropoulos / <i>Digital typography in the new millennium: Flexible documents by a flexible engine</i>
	13	Włodzimierz Bzyl and Tomasz Przechlewski / <i>Migrating to XML: The case of the GUST Bulletin archive</i>
	19	Tomasz Przechlewski / <i>Managing TeX resources with XML topic maps</i>
	26	Luca Padovani / <i>Interactive editing of MathML markup using TeX syntax</i>
	35	Jan Holeček and Petr Sojka / <i>Animations in pdfTeX-generated PDF</i>
	42	Mustapha Eddahibi, Azzeddine Lazrek, and Khalid Sami / <i>Arabic mathematical e-documents</i>
Fonts	48	Mostafa Banouni, Mohamed Elyaakoubi, and Azzeddine Lazrek / <i>Dynamic Arabic mathematical fonts</i>
	54	Karel Piška / <i>Creating Type 1 fonts from METAFONT sources: Comparison of tools, techniques and results</i>
Omega	65	Yannis Haralambous and Gábor Bella / <i>OpenType and Ω: Past, present and future</i>
	68	John Plaice and Paul Swoboda / <i>Moving Ω to a C++-based platform</i>
	71	Jin-Hwan Cho and Haruhiko Okumura / <i>Typesetting CJK languages with Ω</i>
Software & Tools	77	Jean-Michael Hufflen / <i>MIBIBTeX: Beyond LaTeX</i>
	85	Jérôme Laurens / <i>iTeXMac: An integrated TeX environment for Mac OS X</i>
	92	Frank-Rene Schaefer / <i>ŠäferTeX: Source code esthetics for automated typesetting</i>
Graphics	99	L. N. Gonçalves / <i>FEATPOST and a review of 3D METAPOST packages</i>
Critical Editions	105	David Kastrup / <i>The bigfoot bundle for critical editions</i>
Philology	111	Johannis Likos / <i>μονοτόνων: Java-based conversion of monotonic to polytonic Greek</i>
	121	Dimitrios Filippou / <i>Hyphenation patterns for ancient and modern Greek</i>
	127	Pablo Rosell-González / <i>The mayan package and fonts</i>
	131	Manasi Athale and Rahul Athale / <i>Using LaTeX to typeset a Marāṭhī-English dictionary</i>
	135	Jagoba Arias Pérez, Jesús Lázaro and Juan M. Aguirregabiria / <i>Basque: A case study in generalizing LaTeX language support</i>
	140	Péter Szabó / <i>Implementation tricks in the Hungarian babel module</i>
	162	Kenneth R. Beesley / <i>Typesetting the Deseret alphabet with LaTeX and METAFONT</i>
News & Announcements	192	Calendar
	196	TUG 2005 conference announcement
	c3	Practical TeX 2004 conference announcement
TUG Business	193	Institutional members
	194	TUG membership application
Advertisements	195	TeX consulting and production services
	196	<i>Guide to LaTeX</i> , 4 th Edition, by Helmut Kopka and Patrick W. Daly