

# Proceedings of the XVII<sup>th</sup> European T<sub>E</sub>X Conference

JULY 5–8, 2006, DEBRECEN, HUNGARY

*(Preprints)*

Copyright © 2006 Magyar T<sub>E</sub>X Egyesület – MaT<sub>E</sub>X  
<http://www.matexhu.org/>

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 this original English permission notice must be included verbatim.

Preprints designed by István Baranyai.



#### PROGRAMME COMMITTEE

Karl Berry  
Thierry Bouche  
Gyöngyi Bujdosó  
Yannis Haralambous  
Hans Hagen  
Taco Hoekwater  
Bogusław Jackowski  
Jerzy Ludwichowski  
Péter Maczó  
Frank Mittelbach  
Bencze Nagy  
Bernd Raichle  
Péter Szabó  
Ferenc Wettl

#### ORGANISING COMMITTEE

László Baranyai  
Gábor Bella  
Gyöngyi Bujdosó  
Luzia Dietsche  
Klaus Höppner  
Péter Hanák  
Tamás Mihálydeák  
Bence Nagy  
Volker R.W. Schaa  
Péter Szabó  
Maria Jolanta Szelatyńska  
Ferenc Wettl

#### PREPRINTS COMMITTEE

István Baranyai  
Gyöngyi Bujdosó  
Péter Szabó  
Ferenc Wettl

#### SPONSORS

DANTE e.V.  
TUG  
GUTenberg  
Faculty of Computer Science, University of Debrecen  
Panem Kiadó  
Gyöngyi Bujdosó





## TABLE OF CONTENTS

István Bencze et al. -	7	<i>Server side PDF generation based on <math>\text{\LaTeX}</math> templates</i>
Thierry Bouche -	13	<i>A <math>\text{\pdf\LaTeX}</math>-based automated journal production system</i>
Gyöngyi Bujdosó -	19	<i>On a typography based online help on <math>\text{\TeX}</math></i>
Hossam A.H. Fahmy -	23	<i>Typesetting the Qur'an and its specific challenges to <math>\text{\TeX}</math></i>
Katalin Fried -	29	<i>The colorful side of <math>\text{\TeX}</math></i>
Hans Hagen -	33	<i>How to deal with <math>\text{\TeX}</math> in unfriendly situations</i>
Hans Hagen -	35	<i>What tools do <math>\text{\ConTeXt}</math> users get</i>
Yannis Haralambous -	41	<i>New hyphenation techniques in <math>\Omega_2</math></i>
Yannis Haralambous et al. -	47	<i>Open-Belly Surgery in <math>\Omega_2</math></i>
Hartmut Henkel -	55	<i>Making better PDF</i>
Taco Hoekwater et al. -	57	<i><math>\text{\METAPOST}</math> Developments</i>
Taco Hoekwater -	61	<i>Opening up the type</i>
Taco Hoekwater -	63	<i>The making of a (<math>\text{\TeX}</math>) font</i>
Ildikó Miklós -	69	<i>KöMaL CD – the execution</i>
Jean-Michel Hufflen -	71	<i><math>\text{\MLBibTeX}</math> meets <math>\text{\ConTeXt}</math></i>
Bogusław Jackowski -	77	<i>Appendix G illuminated</i>
Siep Kroonenberg -	85	<i>Font installation the shallow way</i>
Siep Kroonenberg -	91	<i>Managing a network <math>\text{\TeX}</math> installation under Windows</i>
László Németh -	97	<i>Automatic discretionary hyphenation in OpenOffice.org</i>
Karel Píška -	99	<i>Font verification and comparison in examples</i>
István Radó -	105	<i>Typography – the art of the letter system</i>
Bernd Raichle -	107	<i>Introduction into <math>\text{\BibTeX}</math> style file Programming</i>
Martin Schröder -	109	<i><math>\text{\pdfTeX}</math> – what was, is and will be</i>
Péter Szabó -	111	<i><math>\text{\LaTeX}</math> programming tutorial</i>
Péter Szabó -	113	<i>Managing a math exercise database with <math>\text{\LaTeX}</math></i>
András Virágvolgyi -	119	<i>Would Aldus Manutius have used <math>\text{\TeX}</math>?</i>



# Server side PDF generation based on L<sup>A</sup>T<sub>E</sub>X templates

ISTVÁN BENCZE, BALÁZS FARK, LÁSZLÓ HATALA, PÉTER JESZENSZKY

University of Debrecen  
Faculty of Informatics  
Egyetem t.  
H-4032, Debrecen, Hungary  
jeszy@inf.unideb.hu

## ABSTRACT

*We present a web-based addressbook application that can generate customized PDF documents using L<sup>A</sup>T<sub>E</sub>X template documents.*

*The application is hosted on a server and users can access its functions using a web browser. (A working L<sup>A</sup>T<sub>E</sub>X system must be installed only on the server side.)*

*Each registered user can manage his or her own addressbook. They can upload L<sup>A</sup>T<sub>E</sub>X templates and can generate multiple PDF documents from a template. Templates are customized to each selected recipient, substituting the appropriate addressbook data element into them. (An example application might be an invitation card or a letter that must be sent to different recipients.)*

*Moreover users can create simple documents (eg. letters) using builtin templates and a simple web-based document editor.*

## INTRODUCTION

The Portable Document Format (PDF) has become one of the most widely used electronic document formats for publishing documents on the Web. It has many advantages that made it very popular. Some of them are the following:

- It is an open standard.
- It is device and platform independent.
- It is suitable for both viewing and printing.
- It is a file format, not a programming language like PostScript. A PostScript file contains code that must be interpreted, whereas a PDF file is rather a description, that results in faster and computationally less expensive processing.
- PDF files are searchable.

The goal of this paper to give an overview of the tools and techniques that can be used to generate PDF documents in Java applications.

The first section presents a brief overview of the family of PDF tools that are available in Java.

In the following section we present our solution that is based on L<sup>A</sup>T<sub>E</sub>X template documents and on the access to an external L<sup>A</sup>T<sub>E</sub>X system.

The last section is devoted to our sample L<sup>A</sup>T<sub>E</sub>X template driven web application that generates PDF documents.

## OVERVIEW OF CREATING PDF DOCUMENTS IN JAVA APPLICATIONS USING CONVENTIONAL TOOLS

This section gives an overview of the widely used solutions for the dynamic creation of PDF documents in Java applications. These tools can be classified as:

- XSL-FO formatters,
- PDF class libraries,
- reporting tools.

In the following we restrict our attention to open source solutions.

### XSL-FO FORMATTERS

XSL-FO is an XML vocabulary for document formatting, a T<sub>E</sub>X-like typesetting language that uses XML syntax. It is a part of XSL, a family of W3C standards for the transformation and formatting of XML documents.

Because of the verbosity of the syntax XML documents using the XSL-FO vocabulary are not edited manually. In order to use XSL-FO one needs an

XML document and an appropriate XSLT stylesheet to transform it into another XML document that uses the XSL-FO markup vocabulary. (The transformation is executed by an XSLT processor, that is commonly available in Java environments.)

Then the XSL-FO document is converted into a readable or printable format by a so called formatter. The most widely used output format is PDF.

Although the current version of the XSL specification became a W3C recommendation in 2001, none of the existing XSL-FO software products (including commercial products) implements the full standard.

Apache FOP [3] is an open source formatter, that is implemented in Java and is a partial implementation of the XSL specification. FOP provides a Java API to access all of its functionality, thus it can be embedded into Java applications without difficulty.

XSL-FO might be a good solution if your data is in XML. There are ready-to-use XSLT stylesheets for standard XML document formats, such as DocBook XML to transform them into XSL-FO. Writing an own stylesheet is not an easy job. Although there are graphical authoring tools a sound knowledge of XSLT and XSL-FO is required.

#### PDF CLASS LIBRARIES

Several Java class libraries are available to create and work with PDF documents. Unfortunately most of them are commercial products.

For example 14 PDF class libraries are listed in the appropriate category of the Google Directory [1] at the present time, and only three of them are available as open source software.

[2] provides a list of open source PDF libraries in Java and contains 6 entries for the time being.

PDF class libraries can be classified as low-level and high level libraries.

Low-level PDF libraries provide low-level access to the contents of PDF documents and allow the creation of PDF documents in Java applications. To work with these APIs the programmer often should be familiar with the PDF document format. It might be very difficult and cumbersome to use them.

High level PDF libraries use object models to model the logical structure of PDF documents. These logical models consist of Java objects that represent building blocks such as pages, chapters and paragraphs. Manipulating the object model programmers can ac-

cess the content of the underlying PDF documents and modify them.

**PJX** A typical example of low-level PDF libraries is PJX [4]. In order to use it one must know all about the PDF document format.

**PDFBOX** PDFBox [5] is a high-level class library. It allows the programmer to access individual pages within a document and to manipulate them. The content of pages can be accessed as a stream of objects, it is easy to add text and images.

Unfortunately the API does not provide access to higher level building blocks such as chapters or paragraphs. For example in order to add some text one must position to the right location within a page.

A major problem of PDFBox is the lack of extensive documentation. The API documentation is quite good, there are also some example programs but there is a need for a developers guide.

**ITEXT** iText [6] is another high level class library that is more user friendly than PDFBox. It uses a higher level abstraction of documents. The building blocks of documents are chapters, sections, paragraphs, list, tables etc. This model looks like a document object model of an XML document. It is well documented, a very good tutorial is also available. According to the website of the project a book on iText will be published by Manning Publications this year.

#### REPORTING TOOLS

These are software tools that can generate business reports based on templates and data in databases and other data sources. Visual report designers may assist in the preparation of the reports. Templates are typically stored as XML documents that can also be edited by hand.

For a comprehensive list of open source reporting tools see [7]. Reporting tools offer varying features and capabilities, for example they support different data sources and output formats. Some of them can produce PDF output and some can not.

**JASPERREPORTS** JasperReports [8] is an excellent and powerful open source reporting tool that is written entirely in Java. It has a Java API that provides full programmatic control over the entire reporting process from report definition to report generation.

Report templates are defined by XML documents

or defined programmatically, but open source and commercial visual report designer tools are available too. Compiled templates can be populated with data that is passed as parameters by the application or that comes from various data sources. A wide range of data sources are supported, such as relational databases (via JDBC and also via Hibernate), EJBs, XML documents and CSV files. When a template is filled the resulting report can be viewed, printed or exported to PDF, XML, HTML, CSV, XLS or RTF.

JasperReport is a professional tool with many other features such as i18n or integrated charting support.

**DATAVISION** DataVision is an open source reporting tool that is very similar to JasperReports. It is also open source and written in Java, and it can be incorporated into a Java application easily. Reports can be created using a visual report designer tool and stored as XML files that can also be edited manually. The generated reports can be viewed, printed and exported to tab- or comma-separated text files, DocBook, HTML, PDF and XML.

Compared to JasperReports it has fewer features, for example supports only relational databases (via JDBC) and plain text files as data sources.

It is mentioned here because to the best of our knowledge it is the only reporting tool that can export to L<sup>A</sup>T<sub>E</sub>X. Note that it uses L<sup>A</sup>T<sub>E</sub>X only as an output format, the user may use the resulting L<sup>A</sup>T<sub>E</sub>X documents to produce PDF or PostScript files. DataVision itself does not interpret L<sup>A</sup>T<sub>E</sub>X files to produce PDF, it uses the iText PDF library instead to generate PDF files directly.

#### PROBLEMS WITH THE ABOVE SOLUTIONS WHEN USING L<sup>A</sup>T<sub>E</sub>X TEMPLATES

There are problems with the solutions that are presented in the preceding section, they are summarized below.

#### PROBLEMS WITH XSL-FO

*Lack of stylesheets in the case of non-standard XML format* If data is stored in a non-standard XML format and a stylesheet is not available to transform it into XSL-FO, it may be a difficult task to create an appropriate stylesheet.

#### PROBLEMS WITH PDF CLASS LIBRARIES

*Lack of flexibility* Documents are created programmatically, any change in the output PDF file requires modification of the source code and the application must be recompiled.

*Difficulty of use* Most low-level PDF class libraries require an in depth knowledge of the PDF format, it may be extremely difficult to generate a PDF file. Even in the case of high-level libraries it may be difficult to achieve the right text layout.

#### PROBLEMS WITH REPORT GENERATORS

*Non-general purpose* They are useful for generating business reports that contain tables and charts, based on data sources. They may not be the best solution to generate conventional documents such as letters. Typesetting large chunks of text and achieving the right layout may be difficult.

#### COMMON PROBLEMS

*Quality* The aesthetic quality of the generated PDF documents are often poor compared with PDF files that are produced by L<sup>A</sup>T<sub>E</sub>X.

#### JAVA-T<sub>E</sub>X INTEGRATION

**ACCESSING T<sub>E</sub>X FROM JAVA** T<sub>E</sub>X and L<sup>A</sup>T<sub>E</sub>X offer the highest typographic quality. They can produce publication-quality PDF files with a professional appearance. It would be very useful if Java applications could benefit from it.

Unfortunately we have no knowledge of any existing standard tool to integrate Java and T<sub>E</sub>X.

Such an integration may work as follows. To produce high quality PDF output a Java application generates a T<sub>E</sub>X file. This is a trivial task since a plain text file must be generated. Then the resulting T<sub>E</sub>X file is passed to a T<sub>E</sub>X system, that will turn it into DVI, PostScript or PDF.

The T<sub>E</sub>X system is accessed by using the `java.lang.Runtime` class, that allows the Java the application to interface with the operating system. The application can have complete control over the T<sub>E</sub>X compilation process, it can interrupt the process if necessary and it also has access to the files that are produced by the T<sub>E</sub>X system.

Extending the above scenario with the use of T<sub>E</sub>X templates offers greater flexibility. In this case the Java application does not generate a T<sub>E</sub>X file from scratch,

but it reads a template and populates it with data. (Report generators also operate in the same way.)

**TEMPLATE ENGINES** According to Wikipedia, a template engine is a piece of software that processes an input text (the template) to produce one or more output texts. Template engines are widely used and very popular in web application development to create dynamic web content. Their most important advantage is that they separate application logic from web page layout.

Many Java-based template engines are available, see [10] for a list of open source Java template engines. They are not used only on the server side to generate HTML, but they may be used in other applications to produce arbitrary textual output, even source code.

## GENERATING L<sup>A</sup>T<sub>E</sub>X SOURCES

**USING FREEMARKER** FreeMarker [11] is a well-known general-purpose open source template engine implemented in Java. Although it is used to generate HTML web pages in servlet-based MVC applications in general, we use it to produce L<sup>A</sup>T<sub>E</sub>X sources based on templates, that are turned into PDF.

FreeMarker has a powerful template language. Directives such as `if`, `switch` and `list` provide programming capabilities, other common programming language constructs as variables, expressions and user defined functions are also available in templates.

Just like in the case of web applications the template engine is frequently used to incorporate database content into templates. In the example below we use Hibernate to access a relational database.

Hibernate [12] is the most popular object-relational mapping (ORM) solution in the Java world. Hibernate provides transparent persistence for Java objects, that allows applications to store, update and delete objects in a relational database. It also provides query and object retrieval facilities. It is simple to use FreeMarker and Hibernate together.

An example is presented here, a L<sup>A</sup>T<sub>E</sub>X table will be produced by FreeMarker from the following template fragment:

```
\begin{tabular}{ll}
\toprule
Title & ISBN\\
\midrule
<\list HibernateUtil.query("from Book b
```

```
where b.year = 2006 order by b.title")
as book>
${book.title} & ${book.isbn}\\
</list>
\bottomrule
\end{tabular}
```

In the example `HibernateUtil` is a helper class whose static `query(String query)` method executes a HQL query<sup>1</sup>, and returns query results as a list of objects. Here we retrieve all books in the database that are published this year sorted by title. The table will contain titles and ISBN numbers of the books and will look like this:

Title	ISBN
Aglaja. Apokrif	9630779668
Kazár szótár	9637448306
Utazás a tizenhatos mélyére	9631425169

**RELATED PROJECTS** Although they have not influenced our work, the NTS and  $\varepsilon_{\mathcal{X}}$ -T<sub>E</sub>X projects must be mentioned here.

NTS stands for New Typesetting System. The goal of the project was to re-implement T<sub>E</sub>X in Java, but it was discontinued.

NTS has been replaced by  $\varepsilon_{\mathcal{X}}$ -T<sub>E</sub>X [13], that is a T<sub>E</sub>X-compatible typesetting system written in Java. Originally it was started as an attempt to enhance  $\varepsilon_{\mathcal{X}}$ -T<sub>E</sub>X, but later the entire system was rewritten from scratch.

The system is under development. Although a downloadable installer is available at the website of the project, the development is currently in pre-alpha stage.<sup>2</sup>

The project is a very promising initiative, but there is much to do. If it will be available it will provide a more flexible Java-T<sub>E</sub>X integration.

The T<sub>E</sub>X Catalogue Online [14] contains two packages that support database access, namely SQL-T<sub>E</sub>X and LaTeXDB.

SQLTeX is a Perl script that reads an input T<sub>E</sub>X file that contains SQL commands and produces an output in which the commands are replaced with the

<sup>1</sup>HQL stands for Hibernate Query Language, it is the fully object-oriented query language of Hibernate.

<sup>2</sup>Namely it is only a development release that is not “feature complete”. The next so-called alpha release will be delivered for software testers for testing.

results. LaTeXDB is a similar preprocessor but it is implemented in Python. Both packages support only MySQL databases.

Input T<sub>E</sub>X files may contain constructs in either case that look like commands (for example `\sqldb`, `\sqlrow` or `\texdbconnection`) but are not T<sub>E</sub>X commands actually. (This means that T<sub>E</sub>X files containing them could not compile.) They will be interpreted and replaced by the preprocessor to produce T<sub>E</sub>X files that should compile without any errors. In that sense SQLTeX and LaTeXDB operates in the same way as FreeMarker, but they use T<sub>E</sub>X syntax.

### A SAMPLE WEB APPLICATION THAT GENERATES PDF

We have developed a web application to demonstrate the above approach in practice. Using the web application requires registration. Each registered user can manage his or her own addressbook and can generate PDF files based on L<sup>A</sup>T<sub>E</sub>X templates. The user selects entries of the addressbook and these are used to populate the template with data. A separate PDF file is generated for each selected entry that will be offered for download in a single ZIP file.

For example this can be used to generate PDF letters that are customized to each recipient. The PDF files are produced by L<sup>A</sup>T<sub>E</sub>X that is a guarantee for quality. Many users do not have L<sup>A</sup>T<sub>E</sub>X installed on their computer, using the web application they get access to a L<sup>A</sup>T<sub>E</sub>X system. (It is also possible to use the addressbook not at all and to produce single PDF files only.)

After logging in users have the following options:

- manage addressbook (add, delete and modify entries),
- upload an existing template and generate PDF(s),
- create a new template with a simple web-based editor and generate PDF file(s).

If the third option is selected the user is presented with a list of predefined templates. These templates are L<sup>A</sup>T<sub>E</sub>X document skeletons that are stored on the computer hosting the web application. The following templates are installed by default: article, book, report, letter, empty.<sup>3</sup> The document editor is initialized with the selected template.

---

<sup>3</sup>Additional templates can be added easily.

The templates that are uploaded or edited by the user should be valid L<sup>A</sup>T<sub>E</sub>X documents that should compile without any errors, although they may contain constructs that have special meaning. Text surrounded by ‘@’ characters is a variable reference, and a replacement text will be substituted for it.

There are variable references that have predefined meaning, for example

- `@current.name@` means the full name of a person in an addressbook entry;
- `@current.name.firstname@` is the first name of a person in an addressbook entry;
- `@current.addresses.country@` is the country of the default postal address of a person in an addressbook entry;
- `@current.addresses.home.zipcode@` means the zip code of the home address of a person in an addressbook entry;
- `@current.phonenumbers.office@` is the office phone number of a person in an addressbook entry.

These variable references can be used to generate multiple PDF files from a single template based on addressbook entries. Any other variable references such as `@signature@` are called static variable references, that will be replaced with static replacement text.

The user is presented with a list that contains all static variable references that occur in the template. For each of them a replacement text may be specified.

The next step is to select the output format, the possible choices are DVI, PostScript and PDF.

If the template does not contain any variable references or contains only static variable references then a single result file will be generated. Otherwise on the last step the user must select at least one addressbook entry, and a DVI, PostScript or PDF file will be generated for each of them.

The results are offered for download in a ZIP file, that contains the generated DVI, PostScript or PDF file(s) together with the log file(s) and L<sup>A</sup>T<sub>E</sub>X source(s).

The `\write18{command}` construct allows the execution of operating system commands and it is a potential security risk. For security reasons the `\write18` feature should be disabled (normally this is the default in T<sub>E</sub>X systems).

The following technologies and software products was used in the development: JDK 5.0, Apache Tomcat, JavaServer Pages (JSP), PostgreSQL, Hibernate. Note that we did not use FreeMarker, there was no need for such a complex template engine in this application.

#### REFERENCES

- [1] A list of PDF class libraries for Java  
[http://www.google.com/Top/Computers/Programming/Languages/Java/Class\\_Libraries/Data\\_Formats/PDF/](http://www.google.com/Top/Computers/Programming/Languages/Java/Class_Libraries/Data_Formats/PDF/)
- [2] Open Source PDF Libraries in Java  
<http://java-source.net/open-source/pdf-libraries>
- [3] Apache FOP  
<http://xmlgraphics.apache.org/fop/>
- [4] PJX  
<http://www.etymon.com/epub.html>
- [5] PDFBox – Java PDF Library  
<http://www.pdfbox.org/>
- [6] iText, a Free Java-PDF Library  
<http://www.lowagie.com/iText/>
- [7] Open Source Charting & Reporting Tools in Java  
<http://java-source.net/open-source/charting-and-reporting>
- [8] JasperReports  
<http://jasperreports.sourceforge.net/>
- [9] DataVision  
<http://datavision.sourceforge.net/>
- [10] Open Source Template Engines in Java  
<http://java-source.net/open-source/template-engines>
- [11] FreeMarker  
<http://freemarker.sourceforge.net/>
- [12] Hibernate <http://www.hibernate.org/>
- [13]  $\epsilon\mathcal{X}$ -T<sub>E</sub>X <http://www.extex.org/>
- [14] The T<sub>E</sub>X Catalogue OnLine  
<http://texcatalogue.sarovar.org/>



# *A pdf $\text{\LaTeX}$ -based automated journal production system*

THIERRY BOUCHE

Cellule MathDoc

UMS 5638 (Université Joseph Fourier & CNRS)

100, rue des Maths

Domaine Universitaire

38402 St-Martin-d'Hères (France)

thierry.bouche@ujf-grenoble.fr

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

## ABSTRACT

*We present the recent development of a production system for mathematical serials with both an electronic and paper version. The challenges were many: (i) no house style layout should be imposed, as the journals come from different publishing houses and may have very different typographical options; (ii) produce screen-optimised and printer-friendly output at once; (iii) avoid any duplication of information so that every aspect of the publications are always in sync (Web site metadata, table of contents...), thus (iv) generate on the fly article's page numbers, XML metadata at the published volume level from one master  $\text{\LaTeX}$  source file tree. Using available technology (pdf $\text{\LaTeX}$ , pdfpages.sty and \write18), the proposed solution to these problems appeared amazingly simple and easy to use. However, we'll show that there is quite some room left for improvement.*

## INTRODUCTION

At the end of fall 2003, discussions began in the French mathematical community about a federated effort for high-quality online publishing of our academically (meaning: independent and learned society) published research journals. One driving force of this project was the achievements of the NUMDAM digitisation program, which has more or less settled standards for delivery and navigation of a significant part of the mathematical literature (see <http://www.numdam.org>).

Among the prominent features of NUMDAM, we have the rich set of metadata for each article, including tagged bibliographies, and the powerful search engine associated to it, written by Claude Goutorbe of Cellule MathDoc. Thanks to various matching tools provided by the AMS or developed internally, whatever sensible link can be provided is added to the Web interface, which is something our users enjoy a lot.

It appeared after some investigations that what was almost straightforward to achieve in a retrodigitisation process could become a nightmare to produce in a born digital environment:

1. Although all the journals under consideration were produced with some flavour of  $\text{\TeX}$ , each

one had a specific format with a widely paper-only approach to the publication process.

2. Although all of them had a Web site, none of them had reliable processes to control whether the metadata exposed on this site was consistent with the reality of the paper issues.
3. Although bibliographies are such a routine object in the learned publication business, we could count over 20 ways of “structuring” them in the  $\text{\TeX}$  files.

It turned out that, although we're now in the 21<sup>st</sup> century, the quite old fashioned copy-paste operation (a typical late 20<sup>th</sup> century hobby) was the main procedure on which all these journals relied for the most typical aspect of serials publishing: exposing the same data in many formats and contexts. Just think a minute about the starting page number of an article, which is a rather critical data if you hope to find it somewhere. It will be printed within the article itself (where it is only determined at the last step of production, as it depends on the lengths of all the same volume's articles before it), probably an inner table of contents, possibly a back cover table of contents, presumably a Web table of content, not mentioning an eventual annual index, or third party uses of the

data, as current contents or indexing databases services, that could nowadays be fed through OAI-PMH or RSS feeds...

For instance, let me give the following (anonymous) example: a respected journal once published a paper which, for some obscure (might be scientific!) reason, was ultimately shortened by a paragraph or two in the proof reviewing stage. It was the first paper of the first issue of its year of publication in that journal, and was paginated 1–27 (hard coded in the  $\text{\TeX}$  source) although its final form had 26 pages. The next article was thus paginated starting p. 29 but the printer saved the 2 white pages. All the 2000 printed page numbers in this volume are wrong except the 26 first of them. Last minute changes and copy-paste are the two devils in journal production. A more obvious example: an accepted paper happens to have a serious scientific failure which is found after all the publishing process has been done. The author informs the journal in a hurry that they have to cancel it, of course they do. Now, all pages numbers are wrong for the following articles, go figure where they have already been disseminated!

In a retrodigitisation process, these issues are just annoying, but all you aim at is to create accurate and structured metadata describing an existing paper collection. Moreover, as it is a batch process on a large amount of similar data, high quality can be achieved at a reasonable cost. Production cost and complexity is an issue for our small journals, which heavily rely on a voluntary effort by researchers on their spare time (as well as *Cahiers GUTenberg* which will enter the scene later on).

## GOOD SOLUTIONS TO COMPLEX PROBLEMS ARE SIMPLE

So. How do you produce a journal in such a way that you have detailed, accurate metadata in a versatile format, a powerful Web site with screen optimised versions of the articles, and yet the same paper version?

After some time spent in reviewing existing more or less full answers to this question, mostly based on scripts, heuristics, external programs or auxiliary files, I happened to find one so simple that I think it deserves to be detailed here. In fact, it is so simple that I feel somewhat stupid to expose it within two lines on p. 16 while I spend the rest of the text to discuss the troubles. At the time of writing this, 11 issues

from 3 journals (including the latest *Cahiers GUTenberg*) were made using this tool.

This solution has been made possible rather recently thanks to the collaborative effort of many talented developers, and although I could achieve this because of the power of  $\text{\TeX}$  macro language, I must confess that I never use  $\text{\TeX}$  itself, but engines that understand an extension of  $\text{\TeX}$ 's primitives, yet have a full macro interpreter onboard, namely: Pdf  $\text{\TeX}$  with  $\text{\code{\write18}}$  enabled and (soon) Tralics.

Here is a short description of the system.

## PRINCIPLES

1. Any metadata is input at most once in the system, preferably in the relevant file.
2. Anything that is not deterministically determined by a given file, should stay away from that file.
3. Anything that can be computed, *should* be computed!
4. Do not reinvent the wheel, do not invent exotic formats that no one will master, stay pragmatic but avoid bottle-necks that would impact versatility of future use or quality of the output.

## IMPLICATIONS

1. A journal is a set of volumes, made of issues, made of articles, plus various other material, mostly constant. The journal description belongs to the journal file, the volume description to the volume file, etc. Notice that a page number is essentially a by-product of a completed issue, except the first page of an issue, it should be set nowhere.
2. As it is the de facto standard of math writing, AMS- $\text{\LaTeX}$  was chosen as the input format, with as few extensions as required by the further processing. Bib $\text{\TeX}$  for the bibliographies.

**USER INTERFACE** Of the relevant parts of a journal, I didn't implement the volume level (this would save copying *one* number) but tried to define an *issue*.

**Claim 1** *An issue is entirely determined by*

- *its bibliographic data (journal, year, month, volume, issue),*
- *its first page number,*

- *the ordered list of the articles,*
- *and optional additional material such as advertising, disclaimers, obituaries...*

I write this (and only this) in the issue file:

```
\documentclass[francais,CG,Volume,
  Couverture]{cedram}
\IssueInfo{}{46-47}{avril}{2006}
\SetFirstPage{1}
\SpecialNo{Les fontes (Brest 2003)}
\begin{document}
\makefront
  \includearticle{edito}
  \includearticle{atanasiu}
  \includearticle{ghassan1}
  [...]
  \includearticle{tombeur2}
  \includearticle{devroye}
  \includepub{pubyannis}
\makeback
\end{document}
```

This might look stupid at first glance, but believe me: `\makefront` makes the front matter of the paper volume (including the table of contents), `\includearticle` includes the corresponding article, `\makeback` makes the back matter, etc.

Each article obeys an AMS-L<sup>A</sup>T<sub>E</sub>X structure:

```
\documentclass[CG,francais]{cedram}
\usepackage{x,y}
\title{Formatons les formats de fonte}
  {Formatons\ les formats de fonte !}
\author{\firstname{Luc}
  \lastname{Devroye}}
\address{McGill University,\ etc.}
\thanks{L’auteur...}
```

Assuming that all the articles and other material are in final form (which means that they are in a directory of their own, and that an error-free source master file compiled with pdf<sub>l</sub>atex has been compiled successfully with all cross-references resolved), when you compile (twice) the issue file, it will produce one big PDF which is made of all inner pages of the paper volume, this is sent to the printer. It will also produce the pages of the cover, and an XML file with all the metadata for this volume. In fact, as a side effect, you’ll also find in each article subdirectory a hyperlinked PDF with a first page added, so that everything is ready at once for shipping both the paper and electronic editions of that issue.

## ARCHITECTURE

L<sup>A</sup>T<sub>E</sub>X is a “document preparation system”, it operates at the document level. I’m not convinced by systems that address the above mentioned issues by considering articles as subdocuments of a master document: they require a high level of normalisation of the sources to avoid conflicts (different macros with same names, cross references, etc.), a lot of redefinitions of standard user commands which is very risky when users like shortcuts, and would yield broken links when you provide an article on its own. You can’t expect that mathematicians will obey such strict rules, nor T<sub>E</sub>Xnicians! Thus the relevant document unit in a journal is an article, preferably isolated in a specific directory in order to avoid conflicts with input of figure names. It should be compiled individually and produce a nicely hyperlinked and searchable vector PDF. The metadata relevant to the article are classical: authors’ data, title, abstract, keywords, subject classification, dates, bibliography. The volume, issue, page numbers are not part of the article itself, as it can be moved at any time without affecting its scientific content, thus without edits. Of course, an article is prepared for a journal, so that info should be present in the article file, and determine the layout and many typographical options.

Starting from the following obvious observation: nobody but T<sub>E</sub>X knows how long an article is, when considering its source, I eventually understood that the only reliable solution for setting error-free page numbers was to ask T<sub>E</sub>X to do so. Of course, you could compile a volume with a perl script that would compute things, compile articles, examine the produced PDF to deduce page numbers, modify the articles, recompile, etc. These are heuristics, and will be broken at the first discrepancy between the paper volume and the model volume assumed by the script. In some sense, as long as articles have a “bibliographical” reference, we’re still in the retrodigitisation paradigm when producing an electronic edition: it is the paper model that endows the article with its metadata, so it is by assembling the paper volume that we can deduce the required data to get the final articles. But, more generally, the same applies to purely electronic serial publications: even the table of contents of an incrementally growing online volume is something that is generated as the last step when the latest article is added. Only

flat repositories like arXiv can bypass this question.

## IMPLEMENTATION

ARTICLES As long as articles are concerned, the cedram class is nothing but amsart, with few features added. Namely:

- some extra metadata fields (provision for bilingualism, journal dates...);
- the automatic inclusion of a configuration file at `\begin{document}`;
- a ‘lastpage’ trick;
- the facility to store the literal  $\TeX$  code of the argument of a macro or an environment content and to write it to auxiliary files in various formats by overloading standard macros;
- many hooks added in the presentation code so that all known layout options can be easily implemented;
- some ad hoc definitions for various theorem styles;
- a journal option that loads the journal file defining all constant metadata and make up for that journal;
- some more class options, mostly for compatibility (by default, the class requires hyperref, pdflatex, T1 encoding, LM fonts...).

There is a light version called ‘special’ for things that should look like an article but do not have its full features (editorial statements, e.g.).

When you compile an article at its final stage, it reads a possible configuration file that might override options and provide the needed metadata (issue info, first page), writes out the screen optimised PDF (with a first page added, kind of offprint cover, meant to identify more clearly the article origin when it will travel the Net on its own), a .cdrsom file which contains a  $\TeX$  command providing all the data pertaining to this article that could be used to generate the corresponding TOC line in whatever format, and a .cdrxml file that contains an XML-like snippet with all the metadata for this article.

VOLUMES A volume is made with the same class cedram, with option ‘Volume’, so that all the typographical options are the ones of the journal. There are some specific options to this mode of operation, like ‘Couverture’ which will generate the cover, ‘CouvTires’ the covers for author’s (paper) offprints...

Let me explain what it does line by line, which will show how it works, and why it is so simple and reliable.

```
\documentclass[francais,CG,Volume,
Couverture]{cedram}
```

This sets the volume mode of *Cahiers GUTenberg*, with French hyphenation patterns for the editorial material surrounding articles, and will generate a cover.

```
\IssueInfo{}{46-47}{avril}{2006}
\SetFirstPage{1}
\SpecialNo{Les fontes (Brest 2003)}
```

These set variables : issue number (CG has no volumes), month and year of publication, starting page number of the first article, title of the issue when relevant. These variables will be hold in memory during the whole  $\LaTeX$  run, as well as written to auxiliary files.

```
\begin{document}
\makefront
```

In article mode, many things happen at `\begin{document}`, but not in volume mode, as far as I can tell! `\makefront` could have been automated here. However, this command just selects `\pagestyle{empty}`, and inputs CG-front.tex, which in turn inputs the definition files for the front matter (title page, administrative data, summary). It also inputs a void file that can be populated locally for special occasion issues. The summary is a container constant source file making use of the issue variables and inputting a summary data file with some fixed name which is generated later on (thus the necessity of two runs to complete an issue). In fact, the summary is a ‘special’ item, thus a complete  $\LaTeX$  file which is compiled during the run in a subprocess similar to the ‘article’ case below. The `\makefront` macro ejects all remaining material to be printed, goes to the next odd page, and sets the page counter of the master document to the value given by `\SetFirstPage`.

```
\includearticle{devroye}
```

This is the main operation, but maybe the simplest one. It is so simple that I copy its definition here:

```

\def\includearticle#1{%
  \IncludeArticle[2]{#1/}{#1}%
  \ifx\@empty\articlesXML
    \gdef\articlesXML{#1/#1.cdrxml}%
  \else \g@addto@macro
    \articlesXML{ #1/#1.cdrxml}%
  \fi
  \ifx\@empty\articlesSOM
    \gdef\articlesSOM{#1/#1.cdrsom}%
  \else \g@addto@macro
    \articlesSOM{ #1/#1.cdrsom}%
  \fi
}

```

As you see, `\includearticle` is just a shorthand for the more general `\IncludeArticle` that assumes that the article's master T<sub>E</sub>X file resides in a subdirectory with the same basename. Moreover, it stores in a macro the list of `.cdrxml` and `.cdrsom` files that will be dealt with at the end of the run. Going back few lines in `cedram.cls`, we have:

```

\def\IssueInfo#1#2#3#4{\tkkv=%
  {\ScreenMode\issueinfo{#1}{#2}{#3}{#4}}%
  \issueinfo{#1}{#2}{#3}{#4}}
\let\articlesXML\@empty
\tkvp={\setpage}
\def\pdflatex{%
  pdflatex --shell -interaction=nonstopmode }
\newcommand\IncludeArticle[3][2]{%
  \cleararticlepage
  \immediate\write18{echo '\the\tkvp
    \the\tkvp{\thepage}' > #2#3.cfg}%
  \immediate\write18{cd #2 && \pdflatex #3}%
  \ifcdr@redoBibtex
    \immediate\write18{cd #2 && bibtex #3}%
    \immediate\write18{cd #2 && \pdflatex #3}%
    \immediate\write18{cd #2 && \pdflatex #3}%
  \fi
  \immediate\write18{cd #2 && \pdflatex #3}%
  \includepdf[pages={#1-},noautoscale]
    {#2#3.pdf}}

```

The main article operation is thus the following.

1. The last page is ejected and, depending on the journal style, we go to the next odd page before dealing with the article.
2. The issue info has been stored globally and is written to the auxiliary file for the article, together with the current page number (a traditional `\write` could have been used here as well).
3. Then, the article is recompiled and will use these infos because it reads the just created `.cfg` file at `\begin{document}`. Optionally BibT<sub>E</sub>X and further pdfL<sup>A</sup>T<sub>E</sub>X calls are executed.

4. Finally, the newly generated PDF is included (except, of course, its first page) in the master volume being produced.

The trick here is that you can trust the page counter of the master document: this is the actual paper volume to be printed! Thus you can reasonably be sure that the value of `\thepage` is the correct value for the first page of the next article, which will be included precisely at this page. And this will keep true next time as you input the final PDF of the article right away.

```

\includepub{pubyannis}
\makeback
\end{document}

```

These just to show that you can add advertising, or anything else. The counterpart of `\makefront` is `\makeback`: it includes almost static pages (instructions to authors, subscription info, etc.). In fact, many things happen at `\end{document}`, which is the only place where everything is known about the issue in final form: an XML file is written by processing of the master's and all article's XML snippets, an summary data is generated by concatenation of all article's summary lines, the cover is built by compiling the adequate template.

## METADATA AND FORMAT QUESTIONS

As long as printer and screen PDF files for the Web are considered, the described system has proved to be working quite effectively. But, when you decide to produce versatile metadata from L<sup>A</sup>T<sub>E</sub>X source, you can expect troubles. All typeset material is done by L<sup>A</sup>T<sub>E</sub>X, thus the above mentioned `.cdrsom` files are perfect thanks to the possibility of redefining whatever macro on the fly to have different views on the same data (for instance, one journal has three summaries in it: one in French, with corresponding abstracts, another in English, both at the beginning of the paper volume, and another one set differently on the back cover, where actual titles are used: this is why I store 9 fields in those `.cdrsom` files). Apart from pragmatic reasons I discussed before, there is a fundamental reason to prefer T<sub>E</sub>X source as the master for all metadata: math authors write their papers with T<sub>E</sub>X, and validate their scientific content on the printed result, this is where last minute corrections happen. If you had a full XML process, outputting T<sub>E</sub>X code after automated transformations, the correction process

would be much more difficult to control, and could yield cases where there is simply no way to obtain the wanted physical representation of the article, which is at the time of writing the only long-term reference for the scientific output of the paper.

The first version of the cedram tools assumed a lot of postprocessing of the ‘pseudo’ XML files output by  $\text{\LaTeX}$ . We had the  $\text{\TeX}$  code somewhat sanitised, textual material converted to utf-8 with variable success, and math expressions exposed as GIFs on our Web interface, thanks to latex2html.

My first idea in this respect was to use the kind of trick that is exploited in hyperref to produce properly encoded PDF bookmarks in order to write Unicode files. I was not able to achieve this myself. I also had a look at  $\text{\TeX4Ht}$  which sounds like a good conversion tool from  $\text{\TeX}$  to XML or HTML+GIF as an alternate presentation format. I gave up because of the huge amount of parameters and files to understand before having some output only similar to my expectations.

I am currently experimenting with Tralics [1], which might be the killer application: instead of asking pdflatex to write out a pseudo XML snippet for each article, that will need further processing, it can easily write structured code tweaked for Tralics, where all the data is the literal unexpanded  $\text{\TeX}$  string, leaving all the conversion process from  $\text{\TeX}$  data to Tralics, which is very good at that.

It even parses Bib $\text{\TeX}$  files, but might easily be used as well to structure the bibliography environments! Namely, here is an excerpt from the file generated by the compilation of our example article.

```
\begin{xmlelement}{article}
\xmlbibtexcite {b8}{8}
\boxed{pagedeb}{149}
\boxed{pagefin}{166}
\begin{xmlelement}{auteur}
\boxed{nomcomplet}{\firstname {Luc}
\lastname {Devroye}}
\boxed{prenom}{Luc}
\boxed{nom}{Devroye}
{\killparcode\begin{xmlelement}{adresse}
{McGill University,\ \ etc.}
\end{xmlelement}} \end{xmlelement}
{\killparcode
\begin{xmlattelement}[fr]{titre}%
Formatons\ \ les formats de fonte !
\end{xmlelement}}
\begin{biblio} \bibitem{b8}J.\textsc
{Andr\ 'e} \pointir << Ligatures \&
```

```
informatique >>, \emph{Cahiers
GUTenberg}, \no22, p.\~61--86,
1995. \end{biblio}
\end{xmlelement}
```

After small configuration thanks to the fact that Tralics understands deeply  $\text{\TeX}$  macros, Tralics will generate a wonderful valid XML, with all text converted to Unicode, and maths to MathML. This XML can be exploited at once on our Web sites. The only remaining question is whether the world is ready for MathML. Recent tests show that the quality of the display of MathML expressions in current browsers has drastically increased: it is now similar to  $\text{\TeX}$  in readability (which relies a lot on fine positioning of sub- or superscripts), it enhances considerably accessibility to the math content on the Net. The only remaining difficulty is that, as Tralics is a full  $\text{\TeX}$  interpreter, it cannot generate easily a mixed format sanitising the text strings to well-formed Unicode XML but keeping a verbatim copy of the math formulas in  $\text{\TeX}$ , which might be the most practical for many of our users in the near future...

```
<?xml version='1.0' encoding='iso-8859-1'?>
<article>
<pagedeb>149</pagedeb>
<pagefin>166</pagefin>
<auteur>
<nomcomplet>Luc Devroye</nomcomplet>
<prenom>Luc</prenom>
<nom>Devroye</nom>
<adresse>McGill University,\ \ etc.</adresse>
</auteur>
<titre xml:lang="fr">Formatons les formats
de fonte !</titre>

<biblio type='flat'>
<bib_entry crossref='cite:b8'>
<reference>8</reference>
<bibitemdata>J.\&\xA0;<hi rend='sc'>André
</hi> « Ligatures & informatique »,
<hi rend='it'>Cahiers GUTenberg</hi>,
no&\xA0;22, p.\~61&\x2013;86, 1995.
</bibitemdata>
</bib_entry>
</biblio>
</article>
```

## REFERENCES

- [1] José Grimm, “Tralics, a  $\text{\LaTeX}$  to XML Translator”, *TUGboat*, 24 No. 3 (2003), Proceedings of Euro $\text{\TeX}$  2003.

# On a typography based on-line help on T<sub>E</sub>X\*

GYÖNGYI BUJDOSÓ  
 Faculty of Computer Science  
 University of Debrecen  
 H-4010 Debrecen, P.O.B. 12  
 Hungary  
 bujdoso@inf.unideb.hu

## ABSTRACT

*People using T<sub>E</sub>X often search for on-line information on T<sub>E</sub>X. Although the on-line systems can be found show the syntax of many commands and environments but do not contain typographical recommendations for them. For example, we can find the command `\underline` and its syntax but there is no any hint at that underlining texts is not recommended in documents even if the text is a section title.*

*A project was begun to develop a typography based T<sub>E</sub>X help at our university. This presentation deals with the main features of the system to develop, such as typographical recommendations, T<sub>E</sub>X commands and environment, T<sub>E</sub>X source codes of designed forms and layouts, and the most problematic grammatical rules.*

## INTRODUCTION

People when using T<sub>E</sub>X sometimes need some help about syntaxes of commands. There are some very good on-line (in order to be precise web based) systems on using T<sub>E</sub>X (e.g. [4], [6]) which assist people how to use a command or environment, how to set the indentation of a paragraph, how to italicize a word or how to modify the labels of an enumeration. That is quite useful, save a lot of time, but is it enough? Is it possible to find information on the web on how big should be the indentation, which words should be in italics or bold, what types of labels should be applied?

It is safe to say that it is hard to find information on typographical recommendation on the web and most of the people do not buy books on typography.

An idea follows immediately: An on-line system should be offered to people which contains more than the technical details for word processing.

## BASIC CONSIDERATIONS

In the curricula in many courses on computer science for beginners (also in primary and secondary schools) contain word processing. It means, in general, teaching techniques on how to use menus, dialog boxes and icons. The methodology is similar when teaching T<sub>E</sub>X, for example, for students on mathematics spe-

cializations. This method of teaching word processing results that students can use some functions of a word processor or T<sub>E</sub>X, but they do not take care on the layout of documents. Usually they use the default settings of one (and only one) style file and bold letters for emphasizing texts, sometimes they use headings or centering – and that's all. The layout of their documents is not harmonic and not aesthetic, sometimes quite structureless.

The first idea was to give information to students on typography when teaching the techniques of word processing. This method was not successful enough, because the student thought that these rules and recommendations were without importance. Transposing the focus from techniques to typography ([2]) was more successful, because the first information describing a form was about typography followed by the technical issues. The students were more motivated and attended to layout of their documents.

Currently many students at universities do not have time for attending courses on word processing, have no money for buying books, so they use mostly the web for getting information in this field (as even so about almost everything). Nevertheless the on-line helps on word processing contain only technical information. How to motivate students and, in general, people to take care on layout?

A solution could be an on-line help that rests upon typography ([3]).

\*Research was supported by DIP Kooperációs Kutató Központ, University of Debrecen, project no. GVOP-3.2.2.-2004-07-0021/3.0.

## VERSIONS, CONCEPTIONS AND PROBLEMS

The *first version* of our „on-line help” was rather a textbook on the web than an on-line help. It was a static system of web pages on some parts of the curricula, a set of

- sample pages
  - containing the description of some techniques,
  - showing a recommended layout of the form in question,
  - has to be reproduced.

The system had two main problems: The most important problem was that the sample pages did not contain the whole source code of the special forms, only the syntaxes and detailed descriptions of the commands (because students had to reproduce the pages). Time was when this type of learning was motivating for our students, but the world has been changed, nowadays people prefer ready-to-use things. The second problem was that it was hard to find a command or anything else in the curriculum. It was time to change the conceptions!

The *second version* was a system of web pages. These pages contained the same topics of the first version, plus

- the source codes of many forms.

A new structure was applied: small pages with few pieces of information in them and with many links to the related topics and forms.

There was no problem while there were only two hundred files and one main stream: going through the curriculum via the presentations of typographical objects. Some problems arose when more topics and other streams had to be done to the system. The descriptions had to be sliced and organized in another way.

The result: hundreds of small web pages (containing many redundant information), hundreds of pictures and thousands of links. After a while, it was impossible to handle and update the system. The system was collapsed, before publishing it. So, it's time to change the conception!

## NEW TECHNIQUES

After working a lot on systems which were not good enough, we planed to develop a free system that could be a help to teachers (even if they were not familiar in

programming) in creating and organizing thousands of files of their curricula. This system was designed to be different from the free Learning Content Management Systems (LCMS). (A free LCMS needs system administrator, a programmer for installing, maintaining and updating it.) This project has been cancelled because of various reasons.

A frame that we are able to use can be a free LCMS. Moodle [5] seems the best frame for our purposes, its language support is quite good and it has several agreeable features and useful functions.

Moodle – as the LCMS frames in general – has many advantages such as tests, logs, chats, possibility of tutoring, white boards, searching engine, etc. Nevertheless it has some disadvantages, too. Its worst feature in our case is one of its best properties: the uniformity – that is each window, each dialog boxes, etc. have the same layout determined by the chosen skin. Another problem with it is that it lacks of some functions that could make the system agreeable. It is too rigid and hard to personalize.

## FORMATS

One of the biggest problem when using on-line systems that people need computer and web access for getting information. For people who work always on computers, it does not cause any problem because they have both of them, in general. For other people, who have no note books in their bags and no free (or unlimited) internet access, getting information from an on-line help can be problematic.

A more convenient way would be if the system can be downloaded and installed if needed (in an easy way) onto a local computer.

Visualization is a common learning technique. The system could support this method of study by offering a printable version of the needed parts of the curricula to the users. These pieces of information should be at choice.

## ACCESSIBILITY

Accessibility is an important question in typesetting texts, too. Many disabled people use T<sub>E</sub>X, so we must help them, pre-eminently people with visual impairment. Initiatives can be read, e.g., at [7], [8], [9], or in Hungarian at [1]. Many of them should be adapted to the on-line help.

The minimal requirements that the on-line help should be satisfied are the following:



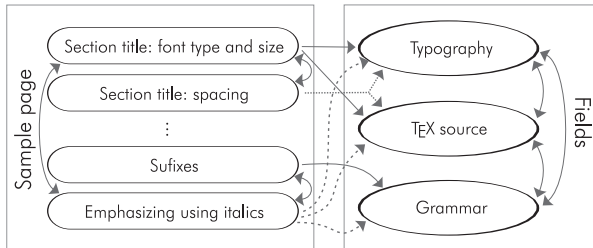


Figure 1: Fields and relations

- It must contain functions for enlarging texts, menu items, icons, pictures, etc. These possibilities could help people with visual and motor impairments.
- Alternative texts should be assigned to each picture.
- It should offer possibilities to users for setting the colors and the contrast applied by the system to the users' wants. This feature is very important, for example, for color- and night-blind people.

#### AUGMENTED CONTENT

Some new topics and many new (cross-)references (see Figure 1) have to be built in. The following new fields are planned to be embedded into the on-line help:

- typographical recommendations,
- commonly (not) used grammar rules,
- a class/style file (or a simple macro file) maker to L<sup>A</sup>T<sub>E</sub>X and plain T<sub>E</sub>X,

and some new features have to be added and lot of relations should be highlighted:

- links from special forms of a displayed page to typographical descriptions and samples that are concerned to the forms (see Figure 2),
- links from special forms to their source codes,
- links to the grammatical rules of problematic words, suffixes, etc.,
- supporting different languages,
- demonstration of designed layouts,
- representing samples on good and bad effectuations of forms and usage of grammar rules. Concerning bad samples of forms must be careful because of legal reasons.

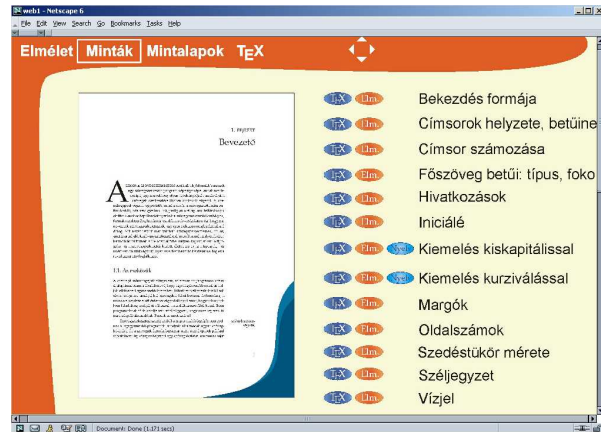


Figure 2: Sample page and relations

#### ABOUT THE PROJECT

The project is consist of two main parts: developing the content, building-up the frame and aligning it to the needs of the content and users.

The content on T<sub>E</sub>X is under development, many of the necessary topics and parts are gathered. The contained descriptions have to be sliced into self-contained pieces taking care on abolishing the redundancies, that is, a system of learning objects has to be developed.

There are many texts on typography, but we try to find a typographer who would write or at least referee them.

The most problematic part is building up and aligning the frame. It needs programmers who have to be supported. Trying to find financial support is in progress.

#### CONCLUSIONS

Our main goal is to develop a system that can motivate people to elaborate nice layouts during their work. Developing an ease-of-use on-line help on word processing combined with typographical issues is a difficult task. It needs a lot of work of several people on various fields: instruction, typography, human-computer interaction, system programming.

In case we develop this on-line system, it will be a good frame to adapt the content to other languages, and offer a place on the web where people can get information not only about T<sub>E</sub>X but typesetting texts and designing nice layouts.

## REFERENCES

- [1] *Accessibility initiatives and principles in Hungarian language*,  
<http://vmek.oszk.hu/vmek2/ajanlas.phtml>,  
<http://www.paramedia.hu/>,  
<http://weblabor.hu/cikkek/iranyelvek>,  
[http://www.w3c.hu/forditasok/wai\\_quick\\_tips.html](http://www.w3c.hu/forditasok/wai_quick_tips.html).
- [2] Gyöngyi Bujdosó, *Teaching word-processing at our university*, Proceedings of Informatika a felsőoktatásban '96 (August 27–30, 1996, Debrecen, Hungary), 1996, pp. 101–109, (in Hungarian),  
<http://www.iif.hu/rendezvenyek/networkshop/96/vegl.html>.
- [3] Gyöngyi Bujdosó, *Online learning course on word processing based on typography*, Proc. of EMES 2005 (May 26–18, 2005, Oradea, Romania), in: *Analele Universităţii din Oradea*, 2005, pp. 33–36.
- [4] W. Macewicz and S. Wawrykiewicz, *Wirtualna T<sub>E</sub>X akademia*, 2004, (in Polish),  
<http://www.ia.pw.edu.pl/~wujek/tex/>.
- [5] *Moodle: A modular object-oriented dynamic learning environment*, <http://moodle.org>.
- [6] *(L<sup>A</sup>)T<sub>E</sub>X Navigator, A (L<sup>A</sup>)T<sub>E</sub>X encyclopaedia*, <http://tex.loria.fr/>.
- [7] *Section 508*, <http://www.section508.gov/>.
- [8] *SENDA: Special Educational Needs and Disability Act*, <http://www.ukcle.ac.uk/directions/issue4/send.html>.
- [9] *WAI: Web Accessibility Initiative*, <http://www.w3.org/WAI/>.

# Typesetting the Qur'an and its specific challenges to the T<sub>E</sub>X family\*

HOSSAM A.H. FAHMY

Electronics and Communications Department,  
Faculty of Engineering, Cairo University, Egypt  
hfahmy@arith.stanford.edu

## 1 PECULIARITY OF ARABIC TYPOGRAPHY

The arabic alphabet has been adopted for use by many languages in Africa and Asia including: Arabic, Dari, Farsi, Jawi, Kashmiri, Pashto, Punjabi, Sindhi, Urdu, and Uyghur. The arabic script is also used for a number of other languages either to present how the language used to be written historically (as for Turkish) or how some write it in an unofficial manner (as for Hausa in western Africa).

Similar to the latin alphabet, with its adoption to several languages, the arabic alphabet acquired new symbols to represent the sounds that do not exist in Arabic. In contrast to the latin alphabet that has dots only on the 'i' and 'j', the arabic alphabet uses dots extensively both above and below the letter shapes to distinguish the different characters. This explicit distinction between the different letters in written text using dots was in itself an addition to the original script which had no dots. In the original arabic script, the distinction between بيت (house) and بنت (girl) when represented as بـ was understood from the context. In general, the developed symbols for the other languages follow the same idea as Arabic and use more dots (up to four) and special marks on the original shapes of the arabic letters.

Arabic being a semitic language, only the consonants are usually written in a word. The equivalent of short vowel sounds are written as additional marks on top of the letters. Obviously, the different languages have different vowels and need different symbols to mark them. In addition to that, since the geographical area covered by the arabic script historically is quite vast, different regions of the world developed different symbols. The result that we have today is a plethora of additional marks developed historically.

A beginner learns that Arabic is written from *right to left* and must practice writing each letter and its connection rules to other letters. Because of the cursive

nature, any letter may connect to the previous and following letters. Hence, a beginner learns the general *four basic forms* of a letter: at the start, middle and end of a word, and isolated. The simplest example of this rule is the equivalent of 'b' in Arabic: بـ. بـ. بـ. بـ. A reader with a sensitive eye would notice that the four shapes of the same letter differ in their width, height above the line, and depth below it. The same structural shape is used for the equivalent of 't' ( تـ ), and 'th' ( ثـ ). The equivalent of 'n' and 'y' share the same shapes as 'b' in the initial and medial forms ( نـ and يـ ) but not in the final or the isolated forms ( نـ and يـ ). In traditional Arabic writing styles (but with the exception of thuluth, riqā', and tawqī' styles [13]), the letters اـ and رـ and their siblings with dots or marks do not connect to the following letter but only to the preceding one.

The cursive nature of arabic script adds another characteristic; many letters combine together to produce new shapes as in الحـ becoming الحـ. In the latin script this phenomena occurs infrequently and when it happens, a *ligature* is used to improve the appearance of the problematic letter combinations such as 'ff' and 'fff'. In arabic typography, on the other hand, the presence of combined letters is abundant *but optional* in many cases. Haralambous [2] gives a long list (yet not exhaustive) of possible 'ligatures' in arabic. While speaking about the history of arabic typography, Milo [9] explains that "each letter can have a different appearance in *any* combination, something that can only be crudely imitated with ligatures". According to Milo [9], most modern books present the connected letter groups "as 'ligatures' and 'artistic expressions' without so much as a hint at traditional morphographic rules". Mackay [8] discusses the range of context evaluation in Arabic and concludes that clusters of four, five, six, and sometimes more letters may combine into a unique shape. Mackay then pro-

\*A graduation project under the author's supervision by: Ibrahim R. Mohamed, Ahmed Z. Ameen, Ahmed M. Amin, Akram A. Mojahed, Kareem O. Sharawi, and Hisham Shihab

poses the use of virtual fonts as an adequate solution.

Another feature in the Arabic script is its reliance on subtle changes to the letter shape to aid the reader in identifying the beginning and end of each letter within a combination. The letter **س** has three “teeth” (vertical pen strokes) similar to the teeth in **ب** and **ت**. When **س** is connected to **ب**, the tail of the **س** may be elongated to alert the reader to the correct grouping of the teeth. Furthermore, the two words **سبع** and **تسع** present different heights for the teeth of the **ب** and **ت** to help the reader as well. This difference in width and height is a type of encoding to prevent a misreading of the word and to aid the trained eye in quickly catching the letter combination. That encoding helps in other cases as well. If due to any reason the dots fade away, a reader faced with **سع** can guess its correct origin. This encoding to emphasize the different letters by raising some teeth is essential in words such as **تسببت** and **تبينت**. However, the raising of the teeth is only possible by investigating the group of letters. So the height of the tooth for **ن** in **تبت** and **تبينت** is not a feature of the individual letter but of the whole combination. The same goes for the height of the dot on top of that same letter as in **سنن** or **سننح**.

In traditional (manual) writing, the “skeleton” of the letter combination is written first then the dots and the other marks are provided. So, a writer probably progresses from the skeleton to the dotted to the vocalized form as **سِنَن** → **سَنَن** → **سُنَن**.

In the Arabic script, a good calligrapher justifies the lines not just by stretching the spaces between the words but mainly by using optional ligatures or wider forms of some letters as in changing **ك** to **كـ** and writing **كتب** instead of **كتب**. The use of optional ligatures and wider forms is the preferred method in high quality works. Another method is to add an elongation to the tail of some letters by using the *taṭwīl* or *kashīdah* symbol ‘**ـ**’ such as **سببـ** instead of **سبب**. This second method has been widely abused in newspapers and low quality materials using mechanical typewriters.

The Arabic script has a large number of writing styles that were developed traditionally to accommodate the different languages and different purposes. Latin scripts use bold, italic, or larger fonts for section headings and for emphasis. Traditional Arabic writings vary the typeface instead. The printings of

the Qur’an as well as of most books almost always use the *naskh* typeface for the main body. The headings, the introductory materials, and the back materials frequently use other typefaces such as the *thuluth*, *ta’līq*, and *ruq’ah*.

In summary, the points just mentioned are:

1. Arabic is written from right to left.
2. Characters in general have four different forms.
3. These forms are of different width, height, depth.
4. The shape (the height of the teeth for example) of a specific form depends on its context.
5. There are additional marks that are put on top or below the character.
6. The horizontal and vertical location of the dots and marks on the characters is not at the same position always but depends on the character *and* its context.
7. There are too many letters that combine (“ligatures”).
8. Ligatures and variable width forms of the letters are used to justify the lines.
9. Several typefaces are needed for special materials in a work of good quality.

With all of these issues, to find a suitable position for the dots and marks on the letter combinations is sometimes a real challenge even for a human let alone a machine.

## 2 AUTOMATED TYPESETTING OF ARABIC

Arabic script does not enjoy the same luxury that Latin script has when it comes to automated typesetting on computers. Milo correctly asserts [9] that the use of individual letters as the building block is not suitable for Arabic. Both Mackay [8] and Milo [10] argue that a layered approach is a better solution. In such a layered approach, some basic elements are provided in the font to represent the skeleton of some letter combinations, an individual letter’s skeleton, or even a part of a letter. These elements are combined first to give the correct skeleton of the word with all the needed shaping for the teeth or other style requirements. On top of that skeleton, a second layer for the dots is added. Then, the vowel marks and any other marks come in

subsequent layers. Milo [11] developed a system with a layered approach for his company, DecoType. It is a proprietary system used by a number of commercial software tools. Due to its proprietary nature, the full details and the extent of the capabilities of this system are not widely known.

Our goal is to provide a freely available system capable of typesetting the Qur'an, other traditional texts, and any publications in the languages using the arabic script. The Qur'an is one of the most demanding arabic texts from a typographical point of view. However, there is a long historical record of excellent quality materials (manuscripts and recent printings) to guide the work on a system to typeset it. Such a system, once complete, can easily typeset any work using the arabic script including those with mixed languages.

Knuth and MacKay [5] were the first to present a working solution for including right to left text (for Arabic and Hebrew) in the T<sub>E</sub>X family. Their proposed T<sub>E</sub>X-X<sub>E</sub>L system is an extension of T<sub>E</sub>X that produces a different DVI file. The enhanced mode of  $\epsilon$ -T<sub>E</sub>X allows bidirectional text processing and produces regular DVI files but  $\epsilon$ -T<sub>E</sub>X does not provide any arabic fonts or any specific functionalities that ease the typesetting of arabic books. Within the T<sub>E</sub>X extensions, both  $\Omega$  [4] and ArabT<sub>E</sub>X [6, 7] have been used for Arabic and have met some of the basic requirements to varying degrees.

With the historical trend to extend T<sub>E</sub>X,  $\Omega$  evolved as an implementation allowing multilingual text processing. As an offshoot of the work on  $\Omega$ , Al-Amal system [2] is an early attempt to typeset the Qur'an specifically. Unfortunately, it is not freely available and its output (as shown in the example published in the paper describing it) falls really short of the desires of a native reader.

Due to its various attractive features,  $\Omega$  was the first choice to achieve our goal. However, the result of our early experiments with the available arabic font provided with  $\Omega$  were not satisfactory. None of the people whom we asked for their opinion liked the font. That font is suitable for a simple publication but definitely not for a high quality work using the arabic script.  $\Omega$  in its current state does not easily lend itself to the layered approach described earlier. The modification of  $\Omega$  is not an easy task since it is a very large system and such a modification means the creation of a new system that is not compatible with

the existing base of T<sub>E</sub>X. The newer developments to  $\Omega$  [3] —once they are stable, widely available, and documented— may help in implementing the layered approach necessary for typesetting high quality texts in arabic.

Lagally in ArabT<sub>E</sub>X [6] preferred to stay within the stable T<sub>E</sub>X standard and perform all the necessary processing with T<sub>E</sub>X macros. That decision allowed ArabT<sub>E</sub>X to be portable to any T<sub>E</sub>X implementation. However, ArabT<sub>E</sub>X had to compromise on the issue of line breaking. For right to left text, ArabT<sub>E</sub>X is forced to handle the line breaking by itself in a slow and complicated algorithm bypassing one of the best parts of T<sub>E</sub>X available for the latin script.

Although not a simple program, ArabT<sub>E</sub>X is confined to a number of style files each performing a specific task. ArabT<sub>E</sub>X implements a layered approach where each character is represented by a skeleton and additional modifiers (dots and vowels). According to the collected opinions, the quality of its font is much better than that of  $\Omega$ . The font still needed improvements but it was an acceptable start. ArabT<sub>E</sub>X uses the L<sup>A</sup>T<sub>E</sub>X license and hence we changed the name of our work to AlQalam (the pen in Arabic).

### 3 IMPLEMENTATION

Our goal of typesetting the Qur'an and traditional texts means a few more challenging requirements in addition to those of the general arabic typesetting. To assist the reader in recitation, several indicators for vowels, joints, text structure, and pausing locations have been added historically to the text of the Qur'an. We present here a few symbols.

*Signs of pause (علامات الوقف) :*

- The ٱ sign : the reader may continue but it is better to pause.
- The ط sign : possible to pause but it is better to continue.
- The waqf jā'iz sign ج : equal possibility to pause or continue.

*Additional diacritics :*

- Ra's khā' ٱ corresponds to sukūn.
- The madda ٱ appears in the Qur'an on many letters such as in كَهَيْجَتِ.
- The small ' in ﴿أَنْ يُشْرَكَ بِهِ﴾ and ' in ﴿وَاللَّهُ عِنْدَهُ حُسْنُ﴾.





Figure 1: The first two lines of surat al-ra'd. An example from the four printed narrations.

Furthermore, there are different “narrations” of the Qur’an that differ in the pronunciation in some locations and hence lead to a plethora of additional marks needed. The vast majority of the printed copies of the Qur’an are in the narration known as Hafs. Only three other narrations (with their special marks for the special pronunciations) are printed in the whole muslim world. The remaining narrations (sixteen remaining for a total of twenty) are still in manuscript form.

Fig. 1 shows an example of the four narrations that exist in print. To make the comparison easier, we present the same two lines from the four narrations written by the same calligrapher and printed by the same press: King Fahd’s complex for printing the Qur’an in Madinah. A simple look at the first word (top right in each example) reveals some of the different symbols needed. Those additional symbols fit well in a layered approach but would be quite difficult to accommodate otherwise. The ج symbol appearing on the first word of the top most narration belongs to the pausing signs. In AlQalam, we introduced a layer for the pausing signs beyond the layer of the dots and the layer of short vowels.

We use the existing symbols in the original font of ArabTeX where they are sufficient. We have improved



Figure 2: Start of surat al-ra'd by AlQalam.

the shape of the madda ‘ـ’ and dagger alif ‘ـ’. The symbols that we have added so far are:



The use of transliteration for a purely arabic document that is a few hundreds of pages long is obviously not practical nor desired. Hence, the default assumption for AlQalam is an input file with the characters coded in Unicode-UTF8. Bi-directional editors such as emacs and gedit (we used both) are good options. Editors have their own limitations though. If a symbol has a unicode point associated with it but there is no key combination mapped to that symbol or no glyph in the editor’s font to represent it, another facility must be used. In a case such as ط the user might supply the unicode as ^db^96. AlQalam interprets this code as belonging to the “signs of pause” category

then raises it to its correct position such as in الْحَقُّ. In addition to that, for some frequently occurring symbols we provide a macro such as \l for ل. The improvement of the input method is one of the major steps in future developments.

Another feature of typesetting the Qur'an is the use of colors. In some printings, certain letters, marks, or sometimes complete words take a different color usually to remind the reader of a pronunciation rule. Educational texts for young children often use similar color encoding schemes to stress new reading concepts and to train their little eyes in picking up the distinctive features of the script. A complete system for dealing with the arabic text should be able to color a piece of a letter combination or some specific marks.

To color the text, we use our modified version of the `acolor` package which Karol Mokry had originally written for ArabT<sub>E</sub>X. With a few modifications to the `aboxes` and `awrite` components of ArabT<sub>E</sub>X, AlQalam allows the user to color the diacritics and the pausing signs of the Qur'an through commands such as `\coldia{blue}`.

#### 4 RESULTS AND FUTURE WORK

Fig. 2 presents the current results of AlQalam for the same narrations shown earlier in Fig. 1.

The initial phase of our work produced a usable system with known problems in the case of some symbols which we are currently solving. In addition to the four narrations of Fig. 2, we completed the work for another ten narrations and wrote a few lines as an example for each one. These results are the first step in a long study of the typesetting requirements for the ten main "readings" of the Qur'an and their expansion into the twenty narrations.

Some of the differences between the narrations constitute simple rules that may be programmed. As an example, every plural pronoun ending with  $\text{م}$  and not followed by  $\text{ا}$  becomes  $\text{م}$  in the reading of 'abī ja'far. In the future, we hope to write macros for all the programmable rules.

Font development is a must. The current results of Fig. 2 are still quite far from the level of Fig. 1. With the vast repertoire of ligatures needed to represent the letter combinations in each typeface (naskh, thuluth, ta'liq, ...), this task is quite complicated. The macros to detect the combinations and produce the appropriate ligatures for the skeletons (which might be changed for line justification) followed by the accumulation of the additional layers of dots and marks is a hard design problem as well.

For line justification, AlQalam inherits the ability to stretch letters using a kashīdah from ArabT<sub>E</sub>X.

However we do not have yet a facility to use the better method of alternative ligatures and wider forms for some letters. In his system, Milo [10] attempts to provide for the use of wider forms. Berry [1] presents another solution (but without high quality fonts and ligatures) to stretch the letters for Arabic, Hebrew, and Persian using `troff`. Within the T<sub>E</sub>X family, Tánh [12] presents what he calls "selective use of multiple glyph" for the latin script although wider forms are not necessarily needed there. Our initial assessment indicates that line breaking for Arabic is an area that needs much more research.

Besides these research problems, AlQalam must extend its capabilities inherited from ArabT<sub>E</sub>X such as footnotes, construction of a table of contents, marginal notes, more L<sup>A</sup>T<sub>E</sub>X commands and environments,...

In conclusion, we provided a summary of the arabic typesetting requirements as well as the first steps of a system to fulfill them. This study is useful for those working on any of the languages written in the arabic script and specifically those supporting the inclusion of Qur'anic quotations. A layered approach is a must for high quality typography. The first version of AlQalam is ready for use but with very limited documentation. Improvements are in process for subsequent versions. AlQalam is free. We would like to share it with anyone interested in testing and providing us with feedback. To the best of our knowledge, there is no other software system that serves the requirements of the different Qur'anic narrations.

#### REFERENCES

- [1] Daniel Berry. Stretching letter and slanted-baseline formatting for Arabic, Hebrew, and Persian with `ditroff/ffortid` and dynamic POSTSCRIPT fonts. *Software—Practice and Experience*, 29(15):1417–1457, 1999.
- [2] Yannis Haralambous. Typesetting the holy Qur'an with T<sub>E</sub>X. In *Multi-lingual computing: Arabic and Roman Script: 3rd International conference — Durham, UK*, December 1992.
- [3] Yannis Haralambous and Gábor Bella. Omega becomes a sign processor. In *EuroT<sub>E</sub>X 2005: Proceedings of the 15<sup>th</sup> Annual Meeting of the European T<sub>E</sub>X Users*, Pont-à-Mousson, France, pages 8–19, March 2005.

- [4] Yannis Haralambous and John Plaice. Multilingual typesetting with  $\Omega$ , a case study: Arabic. In *Proceedings of the International Symposium on Multilingual Information Processing, Tsukuba*, pages 63–80, March 1997.
- [5] Donald E. Knuth and Pierre A. MacKay. Mixing right-to-left texts with left-to-right texts. *TUG-Boat*, 8(1):14–25, 1987.
- [6] Klaus Lagally. ArabTeX — typesetting Arabic with vowels and ligatures. In Jiří Zlatuška, editor, *EuroTeX 92: Proceedings of the 7th European T<sub>E</sub>X Conference*, pages 153–172, Brno, Czechoslovakia, September 1992. Masarykova Universita.
- [7] Klaus Lagally. ArabTeX: a system for typesetting arabic. In *Multi-lingual computing: Arabic and Roman Script: 3rd International conference — Durham, UK*, page 9.4.1, December 1992.
- [8] Pierre A. MacKay. The internationalization of T<sub>E</sub>X with special reference to Arabic. *Proceedings of the IEEE International Conference on Systems, Man and Cybernetics*, pages 481–484, November 1990. IEEE catalog number 90CH2930-6.
- [9] Thomas Milo. Arabic script and typography: A brief historical overview. In John D. Berry, editor, *Language Culture Type: International Type Design in the Age of Unicode*, pages 112–127. Graphis, November 2002.
- [10] Thomas Milo. Authentic arabic: A case study. right-to-left font structure, font design, and typography. *Manuscripta Orientalia*, 8(1):49–61, March 2002.
- [11] Thomas Milo. Ali-baba and the 4.0 unicode characters. *TUGBoat*, 24(3):502–511, 2003.
- [12] Hàn Thê Thánh. Micro-typographic extensions to the T<sub>E</sub>X typesetting system. *TUGBoat*, 21(4):317–434, 2000.
- [13] Mohamed Zakariya. أنماط الحرف العربي. *Al-Computer, Communications and Electronics Magazine* الكمبيوتر والاتصالات والإلكترونيات, 22(8):48–53, October 2005.



# The colourful side of T<sub>E</sub>X

KATALIN FRIED ET AL.

Eötvös Loránd University  
Faculty of Science, Institute of Mathematics  
kfried@cs.elte.hu

## THE DAY IT STARTED...

First of all I would like to say thanks to all who helped me in my understanding and learning T<sub>E</sub>X. (My father, Ervin Fried, my husband, Lehel Juhász, friends, Gabriella Köves, Tamás Bori.)

It all happened in the late eighties. That time Lehel worked as an editor at the publishing house of the Hungarian Academy of Sciences. One day he came home and with a smile on his face he said: I have something you are going to like. (You have to be aware that I started writing books right after I finished university—early eighties—and by this time I had been through two books.) This is, he said to me, a typesetting program. So what, I said, a typewriter can do it. But you can typeset mathematics symbols with it, he said. I can type mathematics on my typewriter, I said. But it looks nice! He added. Now you are talking!, I said, then I simply do not believe you.

## QUESTIONS

Can you type a sum sign?, I asked. Yes, he said. Can you type integrals? Yes. Can you type matrices? Yes. Can you..., can you..., can you..., I kept asking. And the answer was always the same: yes, yes, yes.

After about a couple of hours I remembered something.

Just a few month before that I had problems about denoting arcs in a book. You know what I mean, taking an arc of a circle between the points  $A$  and  $B$  you would like to refer to this arc and denote it in a similar way as you denote a line segment:  $\overline{AB}$  but with an arc above the letters. So I asked, can you typeset such an arc? No, he said. But I can typeset  $\widehat{AB}$  or  $\widetilde{AB}$ . That is not good, I want to see an arc!!! Like  $\widehat{AB}$ . Just much nicer!

But how can you do such a thing?

And there is more! When simplifying fractions I need to cross over the numerator and the denominator. Then I have to write the new numerator and denominator above and below the fraction, respectively.

You know, like

$$\frac{6}{8} = \frac{\overset{3}{\cancel{6}}}{\underset{4}{\cancel{8}}} = \frac{3}{4}.$$

But *how* can you do it?

And can you put a frame around a text? (You must not forget that that time we only had plain T<sub>E</sub>X and no utilities.) You know, like  $\boxed{A \neq B}$ . Or rather  $\boxed{A \neq B}$ .

Yes, yes, but *how do you do it?*

## SOME ANSWERS

That time when these question arose we had no help at all. All we had was a T<sub>E</sub>X book—and only I read English. For framing things we simply had to put them into a box then “wrap them into hrules and vrules”.

```
\def\boxit#1#2#3\hfill\break
{\vbox{\hbox{\vrule\vbox{\hrule%
\kern#2\hbox{\kern#3#1\kern#3}
\kern#2\hrule}\vrule}}}
```

It did not take more than a couple of months to solve this problem. And refining took only another 2-3 years.

Now, crossing over the numerator and denominator of a fraction took somewhat more time. We had to wait until Eberhard Mattes created his version of T<sub>E</sub>X in 1990: emT<sub>E</sub>X.

Its `\special` feature gave us the freedom to create new graphic objects. With these we could solve some of our problems—similar to the framing problem.

We could define nodes:

```
\def\node#1{\special{em:point #1}}
```

We could draw lines between two nodes:

```
\def\line#1#2{\special{em:line #1,#2}}
```

Fractions could be “simplified graphically”. Only we had to measure the numerator and the denominator in T<sub>E</sub>X—with boxing it then measuring the box itself—and draw a line between the appropriate nodes.

And we could even import bitmap graphics into  $\text{\TeX}$ . This was fairly important because we could not draw everything under  $\text{\emTeX}$ .

And from this point on we could create drawings of polygons and lines. We could add letters and symbols to our drawings.

We only had to face one serious problem: how to position the nodes we want to use.

```
\def\put(#1,#2)#3{\vbox to0pt
{\vss\kern#2pt\kern#2pt\hbox
to0pt{\hss\kern#1pt\kern#1pt
\vbox{#3}\hss}\vss}}
```

made us the trick. Notice that this is basically the same idea as the one we used in framing.

But we still could not draw arcs. So we still did not have the arc above  $AB$ .

It was about this time when we discovered a drawing utility for  $\text{\TeX}$ :  $\text{\PCTeX}$ .

It had wonderful features:

1) You could define a plot area to be used. The picture had height, depth, and width contrary to our construction. Why is it important? Because an object that has no height, depth, and width is difficult to input into your  $\text{\TeX}$  file. As soon as you have to position the picture and the text you are going to face problems.

If it has height, depth, and width you can handle it as a  $\text{\TeX}$  object and so you can fit it into your text. True, it is still a hard job. (Of course it is easier if you just centerline the picture.)

2) You could draw circles and ellipses. What a joy! We could do elliptical arcs, circular arcs (parts of ellipses, circles). Alas, no arc above mathematical objects.

And what did we lose when switching to  $\text{\PCTeX}$ ? We had no nodes. That was a great loss so we started to use the two drawing programs at the same time.

$\text{\PCTeX}$  can be used under  $\text{\LaTeX}$ . But we got stuck in plain  $\text{\TeX}$ . Forever, it appears...

## DEMANDS OF AUTHORS

Years have gone and we solved more and more problems concerning drawings. Seeing this, our authors have become greedier and greedier. Not just would have they liked us to create real drawings by computer (ones a graphic artist should do) but also, they would have liked to have colours added to their books.

Luckily, Lehel had his diploma in art, so he could

create all kinds of drawings (only he didn't have time to draw, as he was busy " $\text{\TeX}$ ing". But we gave it a try. There were two things we tried:

1) Drawing, scanning, retouching and importing the drawing into  $\text{\TeX}$ .

2) Drawing by a graphic program and importing the drawing into  $\text{\TeX}$ .

(After these experiments we could import "anything" into  $\text{\TeX}$ .)

We imported bitmap drawings even before. But the age of bitmap graphics was declining. We had to change for eps form. Luckily, we found Tomas Rokicki's epsf.sty file from 1989. (We started to use it many many years later—bitmap did the trick for us for quite a while.)

The drawings were created in some graphic programs (like Illustrator, CorelDraw, etc.).

On the other hand, there was need to do more mathematical objects, such as the graphs of functions and constructions.

We could not keep pace with the demands our authors set for us.  $\text{\emTeX}$  and  $\text{\PCTeX}$  was simply not enough.

But just around that time we found another software perfect for our needs. This software was PSTricks by Timothy Van Zandt (from 1993). For using this we also needed a dvi  $\rightarrow$  ps driver. Tomas Rokicki's "dvips" offered us the postscript output.

What did we gain from it? Everything! It had all sorts of graphics abilities—all of PostTricks. We could embed PostScript codes into the drawings! What a joy we had!

A whole new world opened in front of our eyes with "posttricks" (pstricks).

## MORE ANSWERS (FINE TRICKS)

True, when constructing we had to face a new problem. Euclid's postulates give us the possibility

- to draw a line going through two given points: could be done.
- to open the compasses to a distance of two given points: not simple but could be done somehow.
- to draw a circle with a given length: not simple but could be done.
- to draw the intersection point (if any) of two lines: lacking!

- to draw the intersection points (if any) of two circles: lacking!

Our friend, Tamás Bori gave us a hand. He created a utility with which we could construct the intersection points. We were drawing happily ever after. . .

No, we were not! We found that we could not draw in 3D. I am not a mathematician for nothing. I studied projective geometry. I know how to do the transformation on the 3D coordinates to create such a drawing. So we wrote and used the program. As curves were not given by their coordinates no curves could be drawn. (I have to mention that about a couple of years ago we found a 3D graphing program written by someone else.)

Programming T<sub>E</sub>X reminds me of a conversation I had with a colleague of mine: At university everybody uses T<sub>E</sub>X and when I told him I write programs in T<sub>E</sub>X with variables and calculations and such, he was astonished. ‘Can you really write a program in T<sub>E</sub>X?’ Well, not all of us do it. But for creating an animation I have to have a variable to calculate the number of phases, to calculate the measurement of objects changing, to calculate the shades of colours to

use. Because that’s what animation is.

And we could draw functions dot-by-dot. And we could create animation.

And what is a presentation? Properly animated pages. So, we can create a presentation. As a matter of fact, we have done such a presentation—like this one.

## OPEN QUESTIONS

I want to draw your attention to an important point: these software had been on from the mid 90’s. I have not seen anybody else using it. I believe that we are offered too much and we can take too little. Each of us finds small bits of all the knowledge that had been created in connection with T<sub>E</sub>X. We, ourselves have created utilities, tools for T<sub>E</sub>X we never published. Imagine what a huge amount of knowledge there must be!

Still one question remains: *how can you put an arc above AB?!?*

Finally, I would like to say thanks to all those who posed questions to me to make me think about T<sub>E</sub>X problems (and solve them, most of the cases).



# *How to deal with T<sub>E</sub>X in unfriendly situations*

---

HANS HAGEN  
Pragma ADE, Hasselt  
pragma@wxs.nl

## ABSTRACT

*How to deal with tex in unfriendly situations (multiple trees, potentially conflicting environments, strictly regulated web-services): tools and methods.*



# What tools do ConTeXt users get

HANS HAGEN

Pragma ADE, Hasselt

pragma@wxs.nl

## A BIT OF HISTORY

When we started working on ConTeXt, MS Windows (and before that 4Dos) was our main platform and it still is for development (we use UNIX on the web and file servers and the Mac for fun). So, when ConTeXt got integrated into the TeX distributions we faced the problem of portability. Since one needs auxiliary programs<sup>1</sup> for e.g. sorting an index, we had written TeXutil, and the lack of a commandline handler made us come up with TeXexec. Both were written in Modula but were rewritten in Perl in order to be usable on platforms other than MS Windows. It was easier to maintain a Perl version than to deal with low level platform issues indefinitely.

When our own as well as user demands grew, we wrote more tools and found out that they could best be written in Ruby. In the meantime TeXexec has been rewritten in Ruby, and relevant parts of TeXutil has been merged into it.

## LAUNCHING SCRIPTS

Starting a script on a MS Windows box can be done using a so-called stub, a small program or command file with the same name that locates the similarly named script. On UNIX some shell magic can be used to do the same or one can fall back on a magic preamble (a Bash/Perl mixture) fooling the operating system into locating and spawning the script using the right interpreter. By now, MS Windows has a convenient file association mechanism (but one has to activate it first) while UNIX needs a (nowadays less path sensitive) shebang line and a suffixless copy of the script.

Nevertheless we decided to come up with a less sensitive approach which also gave us the opportunity to accomplish a few more things: TeXMFstart. This script locates and executes a script (or program) in the TeX tree and executes it.

```
texmfstart texexec somefile.tex
```

When you incorporate TeX in workflows, call-

ing TeXexec this way is rather future safe. Actually, because of this method, we could make the transition from TeXexec being a Perl script to being a Ruby program without too much trouble. A side effect of this way of launching scripts is that nested calls are faster because some information is passed on to child runs.

The script is also able to sort out a couple of things, for instance where files reside. Nowadays one will seldom use TeX alone and not all text processing (or related) programs have a clear concept of resource management and/or can work well with a tds conforming tree.

```
texmfstart bin:xsltproc --output=new.xml \
  kpse:how.xsl old.xml
```

This<sup>2</sup> will locate the file how.xsl in the TeX tree and expands the filename to the full path. That way one can keep xslt scripts organized as well. There are a few more such prefixes.

Other features are locating and showing documentation and launching editors with files located in the tree. The following call will open the texmf.cnf file that is currently used.

```
texmfstart --edit kpse:texmf.cnf
```

The script can initialize a tree so one can effectively run multiple trees in parallel. It does so by loading (when present) a file with variable specifications (later more about that).

```
texmfstart --tree=e:/tex-2003 \
  texexec somefile.tex
```

We often use a different tree for each project because commercial fonts may be project related and this way we can move a tree around without running into copyright problems (read: installing all fonts on each box).

```
texmfstart --tree=e:/tex-2003 \
  texexec somefile.tex
```

Another handy feature is conditional processing. In the following case the test file will only be processed

<sup>1</sup>We will use the terms scripts and programs interchangeably.

<sup>2</sup>Please note that a backslash at the end of line denotes a continued line.

when it has changed.

```
texmfstart --ifchanged=test.r --direct R \
  "-q --save --restore < test.r"
```

In a similar fashion one can make running dependent on time stamp comparison. More details can be found in the manual.

## MANAGING CONTEX RUNS

The `TeXexec` script manages a user's `TeX` run. There are many factors that influence such a run:

- Since `ConTeXt` uses the same format for all backends, it depends on loading the relevant backend driver modules. Although one has complete control, life can be made easier when this is done automatically.
- A first pass may generate data needed in a successive pass. There may be references, tables of contents, indices, etc. so we need a way to manage multiple runs. We have to make sure that neither less nor more runs than needed take place.
- A run may demand further action between runs, like graphic manipulations or delayed `MetaPost` execution.
- We may want to run different `TeX` engines, apply different backends, use different user interfaces. Also, the name and way of calling `TeX` may change over time, something that we don't want users to be bothered with.
- We may want to process a `TeX` or `xml` file under different style regime or enable stylespecific modes.
- The document may need an additional page imposition pass, managed in such a way that no auxiliary data gets messed up.
- We may want to close and open the result in a viewer after the run is done.

This and a bit more is handled by `TeXexec`. When dealing with `ConTeXt` files the script will do a few things users are normally not aware of, like making sure that the random seed is frozen for a run, bugs in programs are caught (as long as needed) and that omissions in the `texmf.cnf` settings are compensated for. In addition `TeXexec` provides a few features for combining and manipulating pdf files.

The latest versions of `TeXexec` also support so-called `ctx` files. These are files in `xml` format that describe a process, additional pre- and postprocessing needed, styles and modules to be used etc.<sup>3</sup> This means that one can easily configure projects without the need to tweak source files or editor setups or give explicit commands. Think of situations where an `xml` file (or bunch of files) has to be converted to another variant in order to be processed. `TeXexec` will only do that conversion when needed. In Figure 1 we show the file that is used in the `MathAdore` project.<sup>4</sup> The source file contains `OpenMath` and what we call 'shortcut math' and after normalizing this to `OpenMath` (first conversion) we convert the math to content `MathML` (second conversion).

The source file contains a reference to this `ctx` file and when `TeXexec` is applied to the source file, it will take the appropriate actions. Such a reference looks like:

```
<?ctx-dir job ctxfile ../mathadore.ctx ?>
```

Here "`ctx-dir`" denotes a `ConTeXt` directive.

When dealing with `TeX` files, `TeXexec` will scan the first lines for comments that serve a similar purpose.

## HANDLING THE UTILITY FILE

For a long time `TeXutil` was called from within `TeXexec` to handle the utility file that collects the index entries, tables of contents, references, etc. Nowadays this functionality is integrated in `TeXexec` which is more efficient. We also took the opportunity to enhance the sorting features so that one can mix language specific sorting rules.

The original `TeXutil` is also responsible for some other manipulations, like analyzing graphics. That kind of functionality has been moved to other scripts and more modern ways of dealing with such issues. Because we were in a transition stage to Ruby scripting, it was a good moment to say goodbye to `TeXutil` and concentrate on building a more extensive set of tools.

## THE TOOLS COLLECTION

Instead of expanding `TeXutil`, we decided to spread functionality over multiple scripts. These can be rec-

<sup>3</sup>Although one can use the `ctx` suffix for `ConTeXt` related `TeX` files, this is normally a bad idea.

<sup>4</sup>This project will provide highly interactive math to schools and is conducted in cooperation with the University of Eindhoven.



```

<?xml version='1.0' standalone='yes'?>

<ctx:job>
  <ctx:message>mathadore</ctx:message>
  <ctx:preprocess suffix='prep'>
    <ctx:processors>
      <ctx:processor name='openmath' suffix='om'>texmfstart
        --direct xsltproc
        --output <ctx:value name='new' />
        kpse:x-sm2om.xsl <ctx:value name='old' />
      </ctx:processor>
      <ctx:processor name='mathadore' suffix='prep'>texmfstart
        --direct xsltproc
        --output <ctx:value name='new' />
        kpse:x-openmath.xsl
        <ctx:value name='old' />.om
      </ctx:processor>
    </ctx:processors>
    <ctx:files>
      <ctx:file processor='openmath,mathadore'>v*.xml</ctx:file>
      <ctx:file processor='openmath,mathadore'>h*.xml</ctx:file>
      <ctx:file processor='openmath,mathadore'>openmath*.xml</ctx:file>
    </ctx:files>
  </ctx:preprocess>
  <ctx:process>
    <ctx:resources>
      <ctx:environment>o-m4all.tex</ctx:environment>
    </ctx:resources>
  </ctx:process>
  <ctx:postprocess>
  </ctx:postprocess>
</ctx:job>

```

Figure 1: A ctx file used in the MathAdore project

ognized by their name: they all end with `tools`. If you call them using `TEXMFstart` there is not much opportunity for conflicts with existing tools.

Each tool comes with a manual, so we will not discuss details here.

**CTXTOOLS** This tool provides ConT<sub>E</sub>Xt related features, like generating generic pattern files (so that we are independent), providing editor syntax checking files derived from the generic ConT<sub>E</sub>Xt interface definition (handy for lexers), generating documentation (from the ConT<sub>E</sub>Xt source code), updating ConT<sub>E</sub>Xt (by downloading an archive and regenerating formats), etc.

**RLXTOOLS** The ‘r’ represents resources, normally graphics, the ‘l’ stands for libraries, and the ‘x’ (indeed) for xml. This tool can analyze graphic files and manipulate resources using other programs. For instance it can be used to downsample files at runtime, to handle special color conversion, and to convert graphic

to formats acceptable for T<sub>E</sub>X. By using the runtime converters one can build workflows without the need to rely on additional scripting. There is a dedicated manual on this topic so we will not bore you here with yet another blob of xml.

**XMLTOOLS** You can use this tool to do a simple analysis on an xml file. Another option is to generate a directory listing in xml format. In both cases, the result can be fed into ConT<sub>E</sub>Xt and used in the process. A more obscure option is to generate images from MathML snippets. This script will without doubt include more features in the future.

**PDFTOOLS** This is work in progress. One can for instance roughly analyze pdf files. It also provides a way to manipulate colors in pdf images but that features is now supported in ConT<sub>E</sub>Xt directly. [nice example]

**TEXTTOOLS** Users will seldom need this tool. It can fix things in a tds compliant tree (for instance when the standard has changed), it deals with a few cross platform issues, it can help you to create so called tpm archives (and is meant for ConTeXt module writers) and it can merge updates into your tree.

**MPSTOOLS** In the future this tool will host the now standalone MetaPost to pdf wrapper (mptopdf) as well as the cropper (both are still Perl scripts).

**TMFTOOLS** This script encapsulates some of the functionality of the Ruby based kpsewhich functionality that we use. In the future we may completely move away from the binary because the script is just as fast or faster when it serializes the database. The script can act as a kpsewhich server. The script can also analyze the tree for duplicates.

**RUNTOOLS** Because T<sub>E</sub>X is multiplatform and because we (need to) run services on multiple platforms, we use this script to do things normally done at the console (shell). It just loads the given Ruby scripts with the appropriate library. We also use this tool to generate the ConTeXt distribution.

**EXATOOLS** This is a more obscure tool. It provides some features related to form based style control and web driven T<sub>E</sub>X processing that we use in projects.

**PSTOPDF** This last one is not a collection like the previous tools. It started long ago as a wrapper for Ghostscript. It still provides this function and over the years we've added quite some filtering to it (we just filter the things that Ghostscript fails on or gets confused from). In the meantime we cheat on the name since it also manages the conversion of bitmap images, especially cached downsampling, using ImageMagick as well as conversion from svg to pdf using Inkscape.

**TEXTFONT** This script has been around for a while now and is used to install (commercial) fonts. It generates metric files, map files, and a demo file so that one can see if things went right. ConTeXt does not depend on (ever changing) map file methods and load map files on demand. You can generate map files for dvipdfmx with the previously mentioned ctxtools.

## MORE

There are a few more scripts, like concheck (simple syntax checking) and texsync (synchronizing mini-

mal distributions) but we will not discuss them here.

## INTEGRATION

When setting up multiple T<sub>E</sub>X trees, the trick is in isolating them as well as possible. Because one can never be sure how distributions set things up, we revert to setting environment variables, which will then take precedence over the settings in a regular texmf.cnf file. In the T<sub>E</sub>XMFstart manual you can find more details on how we take care of this, here we only show an example of such an file on Figure 2.

When the tree flag is given, T<sub>E</sub>XMFstart will read this file and set the environment variables accordingly before it launches the program it is supposed to start. In fact, a tree specification can specify a file, but by default the setuptex one is taken.

```
texmfstart \
  --tree=f:/minimal/tex/setuptex.tmf \
  texexec test.tex
```

Since T<sub>E</sub>XMFstart can load multiple such files, we can also use this method to preset more environment variables, for instance pointers to resource like graphics. This is what the -env or -environment line is for, as in:

```
texmfstart --tree=f:/minimal/tex \
  --env=xyz.tmf texexec test.tex
```

The advantage of this variable setting game is that, instead of cooking up scripts with statements like:

```
thread.new do
  ENV["something"] = "nothing"
  a = "texmfstart --tree=f:/minimal/tex --"
  system(a+"env=xyz.tmf texexec test.tex")
end
```

we can put the variable definition in a file and say:

```
thread.new do
  a = "texmfstart --tree=f:/minimal/tex --"
  system(a+"env=xyz.tmf texexec test.tex")
end
```

This has not only a big advantage in terms of isolation (and maintenance) but is also more robust since one can never be sure if another thread is not setting the same variable too, thereby creating confusing (and problems) for all other threads that use the same variable. Since T<sub>E</sub>XMFstart runs as a separate process, it can set its variables independently.

Whenever (on the ConTeXt mailing list) you see mentioning of something named setuptex, you can be sure that it relates to initializing a T<sub>E</sub>X tree (probably a minimal ConTeXt tree) in an isolated way.

```

# file   : setuptex.tmf (the less generic version have suffixes like cmd, sh, csh etc)
# author : Hans Hagen - PRAGMA ADE - Hasselt NL - www.pragma-ade.com
# usage  : texmfstart --tree=f:/minimal/tex ...
#
# this assumes that calling script sets TEXPATH without a trailing
# slash; %VARNAME% expands to the environment variable, $VARNAME
# is left untouched; we also assume that TEXOS is set.

TEXMFMAIN      = %TEXPATH%/texmf
TEXMFLOCAL     = %TEXPATH%/texmf-local
TEXMFFONTS     = %TEXPATH%/texmf-fonts
TEXMFPROJECT   = %TEXPATH%/texmf-project
VARTEXMF       = %TMP%/texmf-var
HOMETEXMF      =

TEXMFOS        = %TEXPATH%/TEXOS%

TEXMFCNF       = %TEXPATH%/texmf{-local,}/web2c
TEXMF          = {$TEXMFOS,$TEXMFPROJECT,$TEXMFFONTS,$TEXMFLOCAL,!!$TEXMFMAIN}
TEXMFDDBS      = $TEXMF

TEXFORMATS     = %TEXMFOS%/web2c/{$engine,}
MPMEMS         = %TEXFORMATS%
TEXPOOL        = %TEXFORMATS%
MPPPOOL        = %TEXPOOL%

PATH           > %TEXMFOS%/bin
PATH           > %TEXMFLOCAL%/scripts/perl/context
PATH           > %TEXMFLOCAL%/scripts/ruby/context

TEXINPUTS      =
MPINPUTS       =
MFINPUTS       =

```

*Figure 2: Example texmf.cnf file*

## CONCLUSION

In this short article we have tried to give you an impression of what is needed in order to make T<sub>E</sub>X usable in a diversity of today's environments. It was not our intention to be complete, because for that purpose we have manuals. One thing should be made clear: although T<sub>E</sub>X itself is pretty stable, the same cannot be said for the environment that it is used in. Just telling T<sub>E</sub>X to process a file is not enough nowadays. This also means that ConT<sub>E</sub>Xt and its tools, in order to keep up, need to be adapted to current needs. On the other hand, by organizing the functionality in tools, and by using a modern and reliable scripting language like Ruby users don't pay a high price for this. Most nasty details can be hidden from them.



# New hyphenation techniques in $\Omega_2$

YANNIS HARALAMBOUS

Département Informatique, ENST Bretagne, CS 83 818, 29 283 BREST Cedex 3, France

yannis.haralambous@enst-bretagne.fr

## ABSTRACT

*By replacing the internal hyphenation engine of T<sub>E</sub>X by an external Omega<sub>2</sub> module, we are able to solve all shortcomings related to hyphenation and to add new features: segmentation of compound words, excentricity, preferential hyphenation.*

## INTRODUCTION

Ever since a computer has hyphenated the word “God” and ruined a night’s sleep of an RCA employee, there has been quite some literature on hyphenation:<sup>1</sup> it is a complex linguistic operation, permanently in use (at least for languages that are hyphenated), and requiring high efficiency. Nevertheless there are some flaws in T<sub>E</sub>X’s approach of hyphenation, as well as some areas where extra features could be helpful.

The flaws are mainly (a) the fact that hyphenation patterns are stored in the format, so that one needs to know in advance which languages will be used in the document and create the appropriate format file, (b) in some contexts, words are not hyphenatable because they are not preceded by glue (for example, a word preceded by a penalty, or the first word of a paragraph), (c) font or color changes prohibit hyphenation, so that a word like “différance” cannot be hyphenated, (d) special hyphenations such as German *backen* becoming *bak-ken* are possible (through the *discretionary* primitive) but cannot be automatic.

New features have been suggested on many occasions: for example, it would be very useful for some languages to prioritize hyphenation between word components rather than between syllables in the same component. In German, the priority list is even three-fold: first comes hyphenation between components, then hyphenation inside the last component, and last *and* least: hyphenation inside the others components. Another useful extra feature is weighted hyphenation: for example, in French, words starting with “con” should not be hyphenated at that location, but this restriction should not be absolute: if one cannot do otherwise, it should be allowed to hyphenate never-

theless (a typical example is the word “conscience” which can be hyphenated only after “con”). So one should be able to specify a penalty value for each hyphen. Another useful feature would be interaction with the user in case of ambiguity requiring morphological, syntactical or even semantic input: a typical case is already stated in the T<sub>E</sub>Xbook: “rec-ord” (the noun) vs. “re-cord” (the verb). Whenever the algorithm detects such an ambiguity, the user should be warned and, why not, asked to disambiguate.

In languages like Greek there is no requirement for hyphenating between word components.<sup>2</sup> Nevertheless, although it is allowed, it looks silly to hyphenate a word such as Παπαχατζηχαλαμπό-πουλος after the two first letters. In such cases it would be useful to prioritize hyphenation towards the middle of the word rather than its borders.

In this paper we present the new hyphenation module of Omega<sub>2</sub>, which solved the problems mentioned and provides infrastructure for the extra features described above.

<sup>2</sup>The reader may wonder why there is such a requirement for German and not for Greek, which uses as least as many compound words as German. Here is a possible explanation: in German, compound words are separated by a glottal stop. For example, *Satzende* will be pronounced “zats[break]ende” to distinguish the components *Satz* (= sentence) and *Ende* (= end). The visually similar *Sitzende* will be pronounced continuously “zitsende” (= sitting). In Greek, however, there is no glottal stop: συνάδελφος (= colleague) is pronounced continuously “sinathelfos” although it is composed by συν (= plus, together) and αδελφος (= brother). This feature of modern Greek language has influenced hyphenation practice. In fact, there are two ways to hyphenate this word in Greek: in modern *démotiké* [5] it will be hyphenated phonetically συν-ά-δε-λ-φος (which contradicts component segmentation), and in *katharevousa* it will be hyphenated etymologically συν-άδε-λ-φος. The question whether *katharevousa* hyphenation patterns should give priority to word components is open.

<sup>1</sup>See, in particular, [6, 7, 8, 1, 10, 11, 12, 13, 14].

## DESCRIPTION

A *module* for Omega<sub>2</sub> is a binary reading and writing horizontal node lists serialized in XML. It is hooked into Omega<sub>2</sub> at two possible locations: one is inside the *end\_graph* procedure (§1096 of [9], just before the call of *line\_break*: that’s when a complete paragraph is sent to the paragraph builder). The second location (which has not been implemented yet) is inside the paragraph builder, for a given vertice of the graph of break nodes.

To say it simply: a module extracts horizontal node lists from Omega’s stomach, modifies them, and puts them back so that typesetting can go on.

But there is something more in Omega<sub>2</sub>: instead of character nodes, we use *texteme* nodes [3], so that we clearly separate glyphs from characters and that we can add arbitrary name/value pairs to each node.

In our case, the hyphenation module will study the paragraph, apply the hyphenation algorithm to textemes and add a “potential hyphenation point” property to some of them. The value of this property is not simply a boolean but rather a penalty, so that the paragraph builder will automatically prioritize some hyphenation points (for example, those between word components).

## PROBLEMS SOLVED

Let us see how our approach solves the five problems described in the first section.

## PROBLEM A: PRELOADING OF PATTERNS

Omega<sub>2</sub> does not preload any patterns in the format file. Patterns contained in external files are dynamically loaded by the module (and subsequently kept in cache), whenever they are needed.

## PROBLEM B: UNHYPHENATABLE WORDS

This problem has always been a nightmare for T<sub>E</sub>X users: now and then, for reasons which seem obscure to the average user, a word will not be hyphenated, resulting in a horrible overfull box. Most of the time the reason is the fact that the word is preceded by something which is not glue. Indeed, according to §891 of [9], a “*potentially hyphenatable part*” of a paragraph is a sequence of nodes  $p_0 p_1 \dots p_m$  where  $p_0$  is a glue node,  $p_1 \dots p_{m-1}$  are character, or ligature or what-sit or implicit kern nodes, and  $p_m$  is a glue or penalty or insertion or adjust or mark or whatsit or explicit kern node (see also [4]).

Since now hyphenation is external to Omega, we can define our own rules. One can apply the original T<sub>E</sub>X rules so that we obtain the same results. Or one can define new rules and obtain potentially better results.

## PROBLEM C: FONT OR COLOR CHANGE

This problem comes from the fact that, as often in T<sub>E</sub>X, the rules for hyphenatable words we described above are not complete. Other restrictions arrive as we read §891: *a whatsit node found after  $p_1$  will be the terminating node  $p_m$  and all character nodes that do not have the same font as the first character node, will be treated as nonletters*. That means that a special primitive or a font change inside a word prohibit hyphenation after them.

In fact, the use of textemes allows us to inject into text properties which will not disable hyphenation. A less typical example is vertical offset: up to now, the T<sub>E</sub>X logo was not a word, but a graphical construction using glyphs. Using texteme properties to specify the lowering of letter ‘E’ it will be at last possible to hyphenate the title of the well-known journal *Die T<sub>E</sub>Xnische Komödie*!

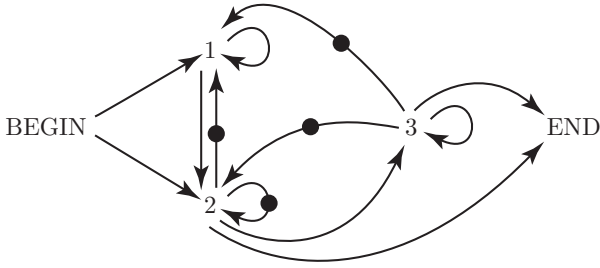
As said in the previous section, since we define our own rules, this restriction can very well be abandoned.

**PROBLEM D: SPECIAL HYPHENATIONS** There is no miracle for that: the various special cases have been hard-coded in the module (one could imagine a syntax for including them in the patterns, but the author considers that their extreme rarity does not justify a syntax enhancement).

## NEW FEATURES

**PENALTIES** As said in [9] §145, a discretionary node produces either a *hyphen\_penalty* or an *ex\_hyphen\_penalty* depending on its pre-break text. This penalty can be changed by the user, but on a global level only, and certainly not separately for each hyphen point. In Omega<sub>2</sub>, there are 65,536 hyphenation penalty registers. Patterns can contain a hyphenation register number (the default register being 0). The hyphenation engine will transmit the highest register number value to the paragraph builder through a texteme property. The paragraph builder will use the penalty value of that register.

It is up to the user to take advantage of these penalties to prioritize hyphenation between word compo-

Figure 1: Finite-state engine of *SiSiSi*

nents, or simply to make the paragraph builder prefer a given hyphenation point rather than some other.

The obvious question which remains is: how do we calculate the various hyphenation penalties in the case of, for example, word component boundaries.

**EXCENTRICITY** To prioritize hyphenation points that occur near the middle of words, we have introduced a number called *excentricity factor*. This number is the slope of a linear function giving the hyphenation register number according to the distance of an hyphenation point from the center of the word: if  $\phi$  is the factor,  $i$  the position of the letter in the word and  $c$  the center of the word, then the hyphenation penalty register will be 0 for letter  $c$ ,  $\text{int}((c \cdot i) \cdot s)$  when  $c > i$  and  $\text{int}((i \cdot c + 1) \cdot s)$  otherwise. So, for example, the word

$\Pi\alpha|\pi\alpha|\chi\alpha|\tau\zeta\eta|\chi\alpha|\rho\alpha|\lambda\alpha\mu|\pi\acute{o}|\pi\omicron\upsilon|\lambda\omicron\varsigma$

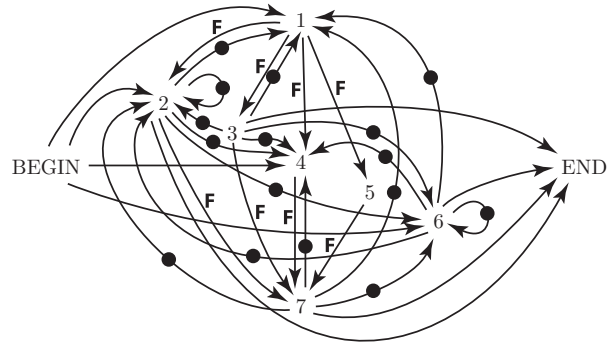
with an excentricity factor of, for example,  $\phi = 0.33$ , will be hyphenated with penalties contained in the following registers:

$\Pi\alpha|_3\pi\alpha|_3\chi\alpha|_2\tau\zeta\eta|_1\chi\alpha|_0\rho\alpha|_0\lambda\alpha\mu|_1\pi\acute{o}|_1\pi\omicron\upsilon|_2\lambda\omicron\varsigma$

It is up to the user to choose the penalty values for each of registers 0, 1, 2, 3. With a higher slope, we get more registers and are able to control hyphenation penalties more finely.

**A FINITE-STATE ENGINE** There are two ways to calculate word component boundaries. Either by using again patterns (as suggested by Antoš in [1]), or by using a finite-state engine, as used by spelling checkers such as *huspell*.

Let us develop the second case. A first approach to word component detection through a finite state engine was *SiSiSi* (= *Sichere sinnentsprechende Silbentrennung für die deutsche Sprache* = “Reliable and Sense-Conveying Hyphenation for the German language”), a  $\text{\TeX}$  extension developped in the early nineties by

Figure 2: Finite-state engine of *huspell*

Barth and Steiner [2] based on their work on hyphenation from the eighties. In *SiSiSi* a word is segmented in three parts: a series of prefixes, a single stem and a series of suffixes. This can be described as the finite-state engine on fig. 1: states 1, 2 and 3 are resp. the one of prefix, stem and suffix. We have twelve transitions  $\text{BEGIN} \rightarrow 1$   $\text{BEGIN} \rightarrow 2$   $1 \rightarrow 1$   $1 \rightarrow 2$   $2 \rightarrow *1$   $2 \rightarrow *2$   $2 \rightarrow 3$   $2 \rightarrow \text{END}$   $3 \rightarrow *1$   $3 \rightarrow *2$   $3 \rightarrow 3$   $3 \rightarrow \text{END}$ , where the asterisk (or • on the figure) means that going through this transition we enter a new word component.

This approach has been proven to be of little efficiency (the *SiSiSi* system was strongly relying on interaction with the user to find and store exceptions to rules).

Another word segmentation approach is the one of spelling checkers. Ispell-based checkers (like *huspell*, which seems to be the most advanced variant) use a “file of affixes” containing lines of the like:

SFX A lig elig [ $\sim$ aeiouhlräüö]lig

which means: there is a set of rules called A containing this rule; this rule says: when you see a word ending by some string verifying the regular expression  $/[\sim\text{aeiouhlräüö}]\text{lig}/$  then strip *lig* and add *elig*. These rules generate new word forms from the existing ones, whenever the latter satisfy the requirements. Similar rules exist for prefixes (starting with PFX).

We have modeled this approach as a finite-state machine with 7 states and 32 transitions. The 7 states are: (1) prefix, (2) stems which are not modified by a SFX or PFX rule, (3) stems which have been modified by a PFX rule (after modification), (4) stems which have been modified by a SFX rule (after modification), (5) stems which have been modified by both a PFX and a SFX rule (after modification), (6) stem which cannot be combined neither by a prefix nor by a suffix,

(7) suffix. It should be noted that only 1, 2, 4 and 6 can be at word begin, and only 2, 3, 6 and 7 can be at word end. The reader can find a graphical representation of this engine on fig. 2.

Here are the transitions: BEGIN→1 BEGIN→2 BEGIN→4 BEGIN→6 1→2F 1→3F 1→4F 1→5F 2→\*1 2→\*2 2→\*4 2→\*6 2→7F 2→END 3→\*1 3→\*2 3→\*4 3→\*6 3→7F 3→END 4→7F 5→7F 6→\*1 6→\*2 6→\*4 6→\*6 6→END 7→\*1 7→\*2 7→\*4 7→\*6 7→END. Those marked by an ‘F’ are conditional transitions: they only apply whenever left and right string belong to the same “family.” Families are necessary because of the regular expressions in SFX and PFX rules.

We will see in the next section how this information is included in the patterns file.

#### PATTERNS FILE

T<sub>E</sub>X users are used to patterns files containing a *patterns* primitive and, in many cases, a *hyphenation* primitive for exceptions. These files sometimes also contain *lccode* commands because only characters with non-zero *lccode* are recognized by T<sub>E</sub>X’s hyphenation algorithm. This part of the code works also as a mapper to equivalence classes, so that patterns can be written using those classes rather than explicit characters. For example, in the case of Greek one can map all combinations of letter *alpha* and diacritics to a single equivalence class and use that class in the patterns. That way, one has lesser patterns and the result is the same (provided, of course, that hyphenation rules are independent of diacritics).

In our new hyphenation patterns files we keep the same pseudo-commands `\patterns` and `\hyphenation`—“pseudo” because these files are not read by T<sub>E</sub>X anymore, but by a tool of our own, called *inittrie*, written in C. These files, for which we recommend the file extension `.pat` are compiled into compressed binary form (file extension `.hyp`) by *inittrie*. This binary form contains a certain number of tries, exactly as formerly stored in T<sub>E</sub>X format files.

Here is a detailed description of the ingredients of `.pat` files.

(LEFT|RIGHT)HYPHENMIN The primitives *left-hyphenmin* and *right-hyphenmin* are used in T<sub>E</sub>X to specify the minimal amount of letters to leave on a line, or to allow on the next line. In our case, these

values are included in the patterns file, so that there is no need to worry about them in the T<sub>E</sub>X file, or babel language file, etc.

EQUIVALENTS The argument of this command consists of a big number of character pairs, separated by blanks. These character pairs play the same role as *lccode* instructions in T<sub>E</sub>X. In each pair, the first character is a character considered as a letter by the hyphenation algorithm, and the second one is its equivalence class. These characters are, of course, all written in Unicode UTF-8. On fig. 3 the reader can see this command displayed under MacOS X.

PATTERNS Patterns are described in the same way as in T<sub>E</sub>X hyphenation files, with an extra convention: numbers between brackets specify the number of the hyphenation penalty register requested. So, for example, the hypothetical pattern

```
.con8s
```

for French language can be replaced by

```
.con8[17]s
```

so that the value of hyphenation penalty register 17 will be used.

SEGMENTHYPHENPENALTY Using this pseudo-command one can specify the hyphenation penalty register to be used between word components. Default value is 1.

SEGMENTHYPHENCHARACTER Through this pseudo-command one can specify the hyphenation character to be used between word components. Default value is the hyphen. In cases like Thai, where we really segment a sentence into words (and words into syllables), the segment hyphenation character will be empty.

SEGMENTPATTERNS The syntax of the argument of this pseudo-command is the same as for *patterns*. These patterns will be used to obtain word components rather than syllables. The other two arguments specify the number of hyphenation penalty register to be used, and the requested secondary hyphenation character.

TRANSITIONS In case we want to use a finite-state engine, we need the commands `\transitions`, `\references` and `\segments` instead of the command `\segmentpatterns`. In the argument of the



[illegible]

Figure 3: Table of equivalents

\transitions command we will write transitions in the syntax given above: strings BEGIN and END for beginnings and ends of words, numbers for all other transitions, the string -> between origin and destination of a transition. Example: BEGIN->1 BEGIN->2 1->1 1->2 2->\*1 2->\*2 2->3 2->END 3->\*1 3->\*2 3->3 3->END for the *SiSiSi* model. The asterisk means that the given transition produces a new segment. A destination followed by an F means that there is a filter: the transition occurs only if origin and destination belong to the same family.

REFERENCES This command contains “references.” A reference is a set of families. The idea is the following: a segment very often belongs to several families (meaning that it can be combined with many other segments, in order to form components and words). Instead of writing all the families to which each segment belongs, we will in fact use a single number. This number will be an index to the set of families to which it belongs, its reference.

The syntax is the following:  
2368=/843/844/845/921/943  
means that segments followed by number 2368 belong  
to families 843, 844, 845, 921 and 943. When checking

whether a transition is allowed, our algorithm will not check if the references are the same, but rather if they have a non-empty common set of families. References are separated by blanks.

**SEGMENTS** Segments are classified by the state to which they belong. The argument of `\segments` looks like:

```

\segments{
1: %prefixes
a2083 abba2084 agyon2084 alá2084 b2085
be2084 bele2084 benn2084 benn-2084
billi62086 c2087 d2088 e2089 egy2086
egybe2084 el2084 ...
2: %original stems
aba3 abafala3 abafalva3 abaffy5 abafi6
abafája3 abajgat8 abakteriális9 ...
3: %pre-altered stems
...
4: %post-altered stems
...
ab2 abafal2 abafalv4 abafáj2 abajga7
abalehot2 abar2 abaújharaszt14
abaújszakal18 abaújszin2 abaújtorn2
abd2 abell25 abelov2 ...
5: %prepost-altered stems
...
6: %stems without affixes

```

```

...
7: %suffixes
abbak2137 abbakat2138 abbakba2138
abbakban2138 abbakb612138 abbakhoz2138
abbakig2138 abbakkal2138 abbakká2138
abbaknak2138 ...
}

```

A number followed by a colon denotes a state. Segments are separated by blanks. They are followed by reference numbers.

Well understood, if the finite-state engine does not require family filtering, then the *references* command will be empty and segments will not be followed by any number.

**LASTSEGMENTPRIORITY** Whenever this option is included in the patterns file, the hyphenation points in the last segment get hyphenation penalty registers increased by a given amount.

**EXCENTRICITY** Gives the excentricity factor.

## REFERENCES

- [1] Antoř D., “PATLIB, Pattern Manipulation Library,” Master Thesis, Masaryk University Brno, 2001. <http://www.fi.muni.cz/~xantos/patlib/thesis/thesis-p.pdf>
- [2] Barth W., Steiner H., “Deutsche Silbentrennung für T<sub>E</sub>X 3.1,” *Die T<sub>E</sub>Xnische Komödie*, 1, 1992, pp. 33-35. [http://www.dante.de/dante/DTK/dtk92\\_1/dtk92\\_1\\_barth\\_steiner\\_deutsche.html](http://www.dante.de/dante/DTK/dtk92_1/dtk92_1_barth_steiner_deutsche.html) and <http://www.ads.tuwien.ac.at/research/\\SiSiSi.html>
- [3] Haralambous Y., Bella G., “Injecting Information into Atomic Units of Text,” in Proceedings of the ACM Symposium on Document Engineering, Bristol, 2005. <http://omega.enstb.org/yannis/pdf/fp10174-haralambous.pdf>
- [4] Haralambous Y., “Voyage au centre de T<sub>E</sub>X : composition, paragraphage, césure,” *Cahiers GUTenberg* 44-45, 2004, pp. 3-53. <http://omega.enstb.org/yannis/pdf/voyage.pdf>
- [5] Haralambous Y., “From Unicode to Typography, a Case Study: the Greek Script,” Proceedings of International Unicode Conference XIV, Boston, 1999, pp. b.10.1–b.10.36. <http://omega.enstb.org/yannis/pdf/boston99.pdf>
- [6] Haralambous Y., “A Small Tutorial on the Multilingual Features of PATGEN2,” in electronic form, available from CTAN as `info/patgen2.tutorial`, 1994.
- [7] Haralambous Y., “Using PATGEN to Create Welsh Patterns,” submitted to *TUGboat*, 1993.
- [8] Haralambous Y., “Hyphenation Patterns for Ancient Greek and Latin,” *TUGboat*, 13 (4), 1992, pp. 457-469. <http://omega.enstb.org/yannis/pdf/ancgreek92.pdf>
- [9] Knuth D.E., *Computers & Typesetting, Vol. B: T<sub>E</sub>X, The Program*, Addison-Wesley, 1986.
- [10] Raichle B., “Hyphenation patterns for words containing explicit hyphens,” CTAN/language/hyphenation/hypht1.tex, 1997.
- [11] Scannell K., “Hyphenation Patterns for Minority Languages,” *TUGboat*, 24 (2), 2003, pp. 236-239. <http://www.tug.org/TUGboat/Articles/tb24-2/tb77scannell.pdf>
- [12] Sojka P., “Hyphenation on Demand,” *TUGboat*, 20 (3), 1999. <http://www.tug.org/TUGboat/Articles/tb20-1/tb62sched.pdf>
- [13] Sojka P., Ševeček P., “Hyphenation in T<sub>E</sub>X — Quo Vadis?” *TUGboat* 16 (3), pp. 280-289, 1995. <http://www.tug.org/TUGboat/Articles/tb16-3/tb48soj1.pdf>
- [14] Sojka P., “Notes on Compound Word Hyphenation in T<sub>E</sub>X,” *TUGboat* 16 (3), pp. 290-297, 1995. <http://www.tug.org/TUGboat/Articles/tb16-3/tb48soj2.pdf>

# Open-Belly Surgery in $\Omega_2$

YANNIS HARALAMBOUS, GÁBOR BELLA, ATIF GULZAR

ENST Bretagne, CS 83 818, 29 283 Brest Cedex 3, France

yannis.haralambous@enst-bretagne.fr, gabor.bella@enst-bretagne.fr, atif.gulzar@gmail.com

## ABSTRACT

*T<sub>E</sub>X and its successors, including the initial version of  $\Omega$ , all suffer from the same technical limitations such as non-adequate support for TrueType/OpenType font formats or the lack of distinction between character and glyph data. In this paper, the authors present  $\Omega_2$  that provides extensibility through both external modules and the texteme concept that supersedes T<sub>E</sub>X's tokens and nodes as well as characters and glyphs.  $\Omega_2$ 's modules, while much more powerful than macros or OTPs, provide relatively easy access to  $\Omega_2$ 's internals without need to touch the source code itself. Among immediate applications are full OpenType support (GSUB, GPOS, etc.), use of independent linguistic tools such as hyphenation algorithms, or support for Unicode's Bidirectional Algorithm.*

## INTRODUCTION

Since its birth, T<sub>E</sub>X has undergone many evolutions that resulted in several extended versions such as  $\epsilon$ T<sub>E</sub>X, pdfT<sub>E</sub>X,  $\Omega_1$ , and others. However, the fundamentals of T<sub>E</sub>X have barely changed: to cite two examples, both its basic text model, that is, the *horizontal node list*, and the concept of the *single main vertical list* are almost exactly the same as 25 years ago.  $\Omega_1$ , as a first step towards Unicode compatibility, introduced 16-bit character codes and some text directionality support but did not change T<sub>E</sub>X's original text model, based on token lists converted to node lists and finally to DVI instructions.

Users and developers have long since recognised the serious limitations of this approach. Inside the belly of T<sub>E</sub>X, character codes included in tokens are replaced by glyph codes, resulting in loss of information if one does not stick to the severely limited T<sub>E</sub>X font encodings, especially in the case of non-Latin scripts: searchability and recovery of the original character data in general become impossible. Support for advanced font formats such as OpenType, essential for writing systems having contextual properties (Arabic, Hebrew, Nastaleeq, the Indic scripts, etc.) but also necessary for some Western typographical features, is also impossible without a clear distinction between the concepts of character and glyph. Still due to the same limitation, until now, no successor of T<sub>E</sub>X could get rid of the TFM font format and provide native support for PostScript- or TrueType-based fonts.

As of today, the most remarkable development

in the T<sub>E</sub>X world regarding support for intelligent font formats is without doubt the XeT<sub>E</sub>X system [1]. However, as far as micro-typography is concerned, XeT<sub>E</sub>X does not have much to do with Knuth's original T<sub>E</sub>X: while the latter is a stand-alone tool, XeT<sub>E</sub>X 'outsources' all the word-level typography to the underlying operating system and external libraries: the ICU library initially developed by IBM [7], ATSUI under MacOS, FreeType under Linux all come into play. So, while XeT<sub>E</sub>X succeeds in combining the OpenType and Apple AAT font technologies with T<sub>E</sub>X's layout and input style, it ties the application to the operating system.

The main reason for preferring such a solution was without any doubt the opportunity to avoid reimplementing the quite complicated Unicode and OpenType engines that had already existed on the operating system level. Moreover, the T<sub>E</sub>X source code itself is far from being easily extendable, despite having been written in the didactic WEB programming language that divides the code into small, easy-to-digest chunks. It lacks modularity and is so highly optimised for performance that the slightest modification can cause a snowball effect of patching and debugging. The single way of extending T<sub>E</sub>X foreseen by Knuth was the creation of new *whatsit* node types, each of which in practice results in a further increase of the programme's complexity.

Unlike XeT<sub>E</sub>X's approach, the developers of  $\Omega_2$  preferred to solve the problem of extensibility by introducing modularity into the system. In fact, OTPs

were already module-like components in  $\Omega_1$ . OTPs are capable of transforming ‘character’ strings into other ‘character’ strings and even of inserting new control sequences. However, due to the early, token-level stage where they intervene, they are limited by the grouping of input text: to show an example, processing of the word ‘*emphasis*’:

`{\it emph}asis`

as a *whole* is not possible, since it is broken by markup into two separate OTP buffers. Even more important is the fact that OTPs are inherently character-level tools and are unable to perform operations other than character substitutions (such as glyph positioning, adding linguistic data, modifying colour, etc.).

Consequently, the objectives that the Omega team has set to themselves were on one hand to solve the character/glyph duality issue by creating data structures capable of storing both, and on the other hand to provide extensibility to  $\Omega_2$  in a more efficient way, in order to allow manipulation of characters, glyphs, and other types of data independently. It will be shown later in the article how these two improvements are tightly related and that they provide the best results when used together.

In the following section, the original text model of T<sub>E</sub>X will be compared to our *texteme*-based approach. Then, external modules will be presented in detail. Finally, it will be shown how using modules together with *texteme* properties opens up possibilities of immediate applications such as OpenType support, linguistic analysis, or fully customised typography beyond the limitations of T<sub>E</sub>X or current font technologies.

## OF CHARACTERS, GLYPHS, TOKENS, NODES, AND DVI INSTRUCTIONS

Text in T<sub>E</sub>X and in its extensions goes through several different states: in the beginning, it is read as *character data* from the input buffer. These characters may either be textual content or T<sub>E</sub>X markup. They will immediately be converted into either *character tokens* or *control sequence tokens*, respectively. A character token consists of the character code and its catcode, while a control sequence is represented by a single identifier. However, the token state is ephemeral: shortly after their creation, both types of tokens are converted

into *nodes*.<sup>1</sup> Character tokens usually become character nodes, but sometimes also ligature nodes. Other types of nodes are also created on-the-fly: kern nodes, glue nodes, discretionary nodes, and so on. Text is organised into horizontal and vertical lists (represented by *hlist* and *vlist* nodes).

An important thing to notice is that fonts come into the picture precisely at the point of converting tokens into nodes. (Ligature, kerning, glue, etc., information all come from font resources.) This is the very moment of T<sub>E</sub>X’s original sin: supposing that

*character code*  $\equiv$  *glyph code from the font*,

characters are being *replaced* by glyphs in character nodes. The reason why hyphenation, usually a character-based operation, still works at a latter stage<sup>2</sup> is this supposed equality that is valid only for a small set of characters, namely for those that were coded in locations common between character and font encodings. Were we to use a different (say, OpenType) font format or a script like Arabic, subsequent character-based operations such as searching, copy-pasting, or hyphenation would all be doomed to failure.

The odyssey is still not over: T<sub>E</sub>X, having done most of its work on node lists, in the end outputs the resulting document using the venerable DVI format where text is encoded through glyph identifiers only, represented on 16 or 8 bits, depending on whether  $\Omega_1$  or another T<sub>E</sub>X-based system is being used. By this time, all the other types of nodes holding non-character data either have already been absorbed in the typesetting process (penalty, discretionary, etc.), or else they now become physical dimensions (kerning, offsets, glue, etc.) or special DVI instructions (specials, etc.). This is the end of the story: our output DVI document is purely presentation-oriented and in no way is able to provide the original character data, long lost in the process.

## TEXTEMES

*Textemes*<sup>3</sup> are one solution to the problems presented above. The idea is to replace T<sub>E</sub>X’s various data representations, namely characters, character tokens, some types of nodes, as well as glyphs in the output, by a single entity: the *texteme*.

<sup>1</sup>OTPs extend the lives of tokens somewhat: they read character tokens and output both character and control sequence tokens.

<sup>2</sup>Namely, at line breaking.

<sup>3</sup>Introduced as *sign* at the EuroT<sub>E</sub>X 2005 conference.

A texteme, as presented in detail in [5] and in [4], is a *set of properties*, where a property is basically a *key-value pair*. A texteme usually represents a character, its glyph, and other related data. More specifically, character code, glyph and font identifiers, and any kind of information related to the atomic units of electronic text are all represented by texteme properties.

How do textemes work in  $\Omega_2$ ? The general idea is to let information accumulate inside textemes instead of converting data from one form to the other. Raw, unformatted text is a texteme string where textemes contain only *character* properties. When raw text (with markup) is fed into  $\Omega_2$ , a *catcode* property is added to every texteme. Once font information were read, instead of creating character nodes, the same textemes—texteme nodes—are carried on, only with new *glyph* and *font* properties added. No ligature nodes are created: an ‘fi’ ligature is represented by two textemes linked together, the first with *character=f* and *glyph=fi*, and the second with *character=i* and *glyph=∅* (empty). Some other information like kerning, glue, or penalties, become texteme properties all the same, resulting in simplified text structure.

Separating character and glyph codes while having access to both of them throughout the whole typesetting process proves to be very useful for tasks such as hyphenation (as it is shown in Yannis Haralambous’s article [3]). However, the ultimate goal is to be able to produce final documents that keep all these accumulated (and useful) information. A document that displays glyphs but also holds the original character-based text has an enormous technical advantage compared to glyph-only documents where retrieval of characters is only possible through non-standard, error prone glyph naming schemes.

The PDF format makes such a double encoding of text possible through the `ActualText` operator: for every glyph or glyph sequence, the corresponding character or character sequence may be defined. This way, even cases like multiple glyphs corresponding to a single character or reordered glyph sequences can be handled correctly, something that would never be possible through glyph naming.

Unfortunately,  $\Omega_2$  does not (yet) produce PDF directly and the DVI format does not offer mechanisms similar to PDF’s `ActualText`. It is therefore not possible to output texteme data into DVI without breaking compatibility with the original DVI format. As a

temporary solution,  $\Omega_2$ ’s DVI format has been slightly modified in order to include texteme-related information that is interpreted by a patched `dvipdfmx` utility that produces PDF documents with `ActualText` operators. This is a quick and dirty solution, but it works.

The document creator is by all means allowed and encouraged to invent and use their own properties in their documents. First of all, the set of available texteme properties is open and extensible.<sup>4</sup> Such user-defined properties can be added either automatically, by linguistic analysers and various text processor tools, or manually, by a texteme-compatible text editor (a simple prototype of which has been developed by students of ENST Bretagne). In this editor, texteme properties are added and manipulated in an intuitive, graphical way. Texteme-based documents can then be saved in XML that  $\Omega_2$  will be able to interpret and thus rebuild texteme strings. (The XML reading capability of  $\Omega_2$  has not been developed yet.)

How do texteme properties come into play during text processing? External modules are the answer.

## EXTERNAL MODULES

### THEORETICAL CONSIDERATIONS

As mentioned before,  $\Omega_1$ ’s OTPs are basically character processors at the token level. This approach is not sufficient when non-character data need to be processed (e.g., glyph substitution or glyph positioning). First, with the introduction of textemes, access to individual texteme properties as opposed to mere character strings becomes necessary. Secondly, even if an extended OTP syntax and input/output scheme allowed the handling of such information, OTPs are still called at the token level where font data have not yet been read by  $\Omega_2$ .

Consequently, in order to allow OTP-like external modules to process font-dependent information, their point of activation needs to be displaced to a later point, to the node level.

But there is a problem: since nodes (including texteme nodes) and node lists are considerably more complicated data structures than characters or tokens,  $\Omega_1$ ’s internal OTP approach is not powerful enough to describe transformations on them. For these reasons, our new external modules need to be standalone binaries that communicate with  $\Omega_2$  using a well-defined

<sup>4</sup>Namespaces are used for semantic disambiguation.

XML format (more on this later). Since these binaries can be written in any programming language, there is no limitation to their computing power, unlike former internal OTPs that were equivalent to finite state automata.

In reality, there are no less than three well-defined points during  $\Omega_2$ 's typesetting process where external modules may be called:

1. on token-based input text (as in the case of OTPs);
2. on yet unbroken horizontal node lists that represent whole paragraphs;
3. on node lists representing individual lines during paragraph breaking.

Each of these three legal intervention points corresponds to a set of well-defined processing tasks. The first point is used by character-level transformers and analysers. They receive a simple list of textemes, uninterrupted by control sequences on grouping braces (no wonder: these tokens act as boundaries for the OTP buffer), and containing mostly character information. They are supposed either to perform character transformations (e.g., converting from a local transcription scheme or encoding to Unicode, preprocessing, etc.) or to generate new texteme properties. (At the moment of writing the article,  $\Omega_2$  is not yet capable of adding texteme properties at this stage.)

The second type of module reads entire paragraphs and operates on node lists. This type of module applies, among others, OpenType glyph substitution and positioning rules. However, as a result of working with nodes instead of just characters, such modules have enormous power as well as responsibility over the behaviour of  $\Omega_2$ : they have full access to every aspect of the text including horizontal and vertical lists, kerning, penalties, and so on. Were a module to, say, substitute a glyph by another, it would have the responsibility to update the corresponding kerning information or at least make sure that this will be done subsequently either by  $\Omega_2$  or by another module.

Finally, the third type of module is called on individual lines, inside the line breaking algorithm. Similarly to modules of the second type, it operates on node lists. Its task is to perform line-related operations such as optical kerning, OpenType JSTF (justification) support, or line-dependent glyph substitutions and positionings (e.g., an OpenType contextual liga-

ture is invalidated by a nearby line break).

The fragility of node lists when manipulated externally may seem worrying. Indeed, it is very easy to produce typographically unacceptable documents and even to freeze  $\Omega_2$  through erroneous or malicious node operations. Creators of modules should respect rules regarding what and in what order they are allowed to modify. Correct ordering of modules is of crucial importance: for example, the order *character transformations* – *glyph substitutions* – *glyph positionings* should always be respected, otherwise regression problems may arise.

Indeed, node-level modules represent a drastic surgical intervention in  $\Omega_2$ 's digestive system: it is as if  $\Omega_2$ 's stomach and intestines were piped into external digestion machines. The reader will kindly excuse the authors for the somewhat disturbing analogy and read on to see how in practice modules are called from  $\Omega_2$ .

**MODULES IN PRACTICE** Module support in  $\Omega_2$  is currently in prototype stage, that is, developer- and user-friendly macros and libraries are only scarcely available at the moment.

External modules are implemented as standalone binaries and communicate with  $\Omega_2$  through signals and an input-output buffer. For performance reasons, these binaries run as *demons*, that is, they are spawned only once and then remain idle until they receive data to process as well as a wakeup signal. Once a module has finished processing, after writing its output into the same input-output buffer, it goes back to sleep again, and  $\Omega_2$  continues by waking up the following module.

More precisely: a module binary is launched by the `\registermodule` primitive. By writing `\registermodule{mymodule}{modbin}{par}{10}` the binary programme `modbin` is run, in the future referenced by the name `mymodule`. The `par` parameter means that it is a type 2 (paragraph-level) module and its number in the execution order among modules of the same type is 10. These four parameters are mandatory. There is a fifth, optional parameter (omitted in our example) that sends its argument to the binary when it is launched; this may sometimes be useful for initialisation purposes.

Modules are *asleep* by default. For performance reasons,  $\Omega_2$  does not send any data to sleeping modules. In order to perform their task, modules need to

be both *woken up* and *activated*. They are woken up and sent back to sleep by the `\wakeup{mymodule}` and `\gotosleep{mymodule}` primitives. Note that paragraph- and line-level modules are woken up for *entire paragraphs*, there is no point in trying to wake them up for shorter text segments. *Awake* but *inactive* modules receive and read all data but let them pass through untouched. They activate themselves when they encounter an *activate* special node, inserted by the `\activate{mymodule}` macro. From this point, they process the text until they either read a *deactivate* node or arrive at the end of the buffer. These *activate* and *deactivate* nodes may of course very well appear inside paragraphs, making it possible to activate modules for text segments as small as individual characters (more precisely, textemes).

Text (textemes and other nodes) is transmitted between  $\Omega_2$  and modules in XML format. The full DTD is not provided here for space reasons, only a small but relevant example is given. As it is shown, for paragraph- and line-level modules,  $\Omega_2$ 's *current font table* is also included in the XML buffer since these modules (e.g., an OpenType engine) usually need access to fonts. The font table is then followed by the node list itself: in our case, a *whatsit*, an empty *horizontal list* and two *texteme* nodes. In this simple example, texteme nodes contain only three properties each, namely the character and the corresponding font and glyph ID.

```
<?xml version="1.0"?>
<buffer version="0.1">
  <preamble>
    <fontlist>
      <font id="51" name="ptmr"
        size="1310720"/>
      <font id="52" name="pala"
        size="655360"/>
      ...
    </fontlist>
  </preamble>
  <odelist>
    <!-- whatsit -->
    <wha st="6" intpnl="0" brkpnl="0"
      pardir="0">
      <lbl></lbl><lbr></lbr>
    </wha>
    <!-- hlist -->
    <hls wd="1310720" dp="0" ht="0"
      shift_amount="0" gse="0" gsi="0"
      go="0" dir="0">
    </hls>
```

```
<!-- texteme -->
<t linkl="0" linkr="0">
  <p n="c">74</p> <!-- char: L -->
  <p n="g">12</p> <!-- glyph -->
  <p n="f">52</p> <!-- font -->
</t>
<t linkl="0" linkr="0">
  <p n="c">101</p> <!-- char: o -->
  <p n="g">142</p>
  <p n="f">52</p>
</t>
...
</odelist>
</buffer>
```

A particular advantage of communicating with modules in XML is that certain modules can thus be implemented by very simple XSLT transformations. This way, the task of implementing a module is reduced to the complexity of writing XSLT code.

## APPLICATIONS OF MODULES AND TEXTEMES

**OPENTYPE SUPPORT** Since quite a long time, the Holy Grail of typesetting systems has been to implement robust support for the TrueType and OpenType font formats. Development has been slow on all platforms, due to the investment required on a very wide scale (full Unicode support, internationalisation, availability of actual fonts). No wonder that no  $\text{\TeX}$ -based system, apart from  $\text{\XeTeX}$ , has even come close yet to full OpenType compatibility: without the profound changes in  $\text{\TeX}$ 's text model described in earlier sections of the present article, or at least similar improvements, OpenType support is not fully possible.

As it was already shown in numerous articles such as [1], the main difficulty of developing OpenType-compatible applications lies in the complexity of operations described in OpenType's GSUB (glyph substitution) and GPOS (glyph positioning) tables. In order to achieve this, apart from providing an appropriate text model, typesetting applications also need a powerful OpenType engine. Fortunately though, a lot of effort has already been made in this direction in the free software community, and thus free and cross-platform OpenType libraries are already available: let us mention the *FreeType* ([6]) and *M17N* ([2]) projects that both offer OpenType-related functionalities. The new  $\Omega_2$  system makes use of both of these libraries<sup>5</sup>.

<sup>5</sup>M17N, far from being just a multilingual typesetting engine, is a whole set of libraries aiming at providing complete multilingual

Basically, two aspects of the OpenType format need to be taken care of in  $\Omega_2$ : reading metrics and performing glyph transformations. Most  $\text{\TeX}$ -based systems solve the former issue by falling back to utilities such as `ttf2afm` that convert TrueType metrics into TFM files, the OpenType and the TrueType formats being compatible as far as metrics are concerned. This solution works but is inelegant from the user's point of view since installation and use of TrueType or OpenType fonts require several conversion steps and editing of various configuration files. The authors have thus decided to implement direct access from  $\Omega_2$  to OpenType metrics, without any need for OFM or other files. At the moment of writing the article,  $\Omega_2$  is already capable of reading TrueType metrics directly from the font.<sup>6</sup>

Glyph transformations, on the other hand, are too big a task to implement inside the monolithic  $\Omega_2$  code. Modules are an ideal place for such operations. Both paragraph- and line-level modules are going to be necessary: paragraph-level modules will perform both font-independent (multilingual preprocessing similar to what Microsoft's Uniscribe library does) as well as font-dependent OpenType GSUB and GPOS glyph transformations. Finally, the same OpenType module intervenes once again at the line breaking phase if necessary; also, JSTF support can be implemented at this point.

**HYPHENATION**  $\text{\TeX}$ 's original hyphenation procedure is called in the line breaking phase: at this point, text is converted into lowercase, ligatures are replaced by their original characters, and pattern matching is performed. Consequently, the hyphenation algorithm is an integral part of  $\text{\TeX}$  that is difficult to modify or customise according to the special needs of different languages, apart from language-dependent pattern sets. In  $\Omega_2$ , textemes and modules allow performing hyphenation externally, in a module, opening up the possibility to use more advanced hyphenation algorithms. The general idea is that the external hyphenation module marks potential hyphenation points in words using texteme properties and at the line breaking stage  $\Omega_2$  simply selects from the marked hyphenation points the ones giving the least badness.

support.

<sup>6</sup>Kerning data are not read by  $\Omega_2$  as this operation is planned to be implemented inside an external module.

See [3] (in this same proceedings volume) for a more detailed description of new hyphenation techniques used in  $\Omega_2$ .

### GETTING RID OF (SOME) $\text{\TeX}$ MARKUP

Through properties, textemes provide a new way of enriching electronic text. In some cases, such properties can substitute markup that would otherwise serve the same purpose. The interest in doing so lies in the simplification of input text, an important gain both from a technical and a usability point of view. Consider the following example of  $\text{\LaTeX}$  code:

The `\verb#\textcolor#` command's purpose is to colorify text, such as this word in `\textcolor{red}{red}`.

There are several problems with the above snippet: first of all, it is far from being intuitive. To typeset the `\` character, the user needs to use the `\verb` command, which is one way of escaping the special role of the backslash character. Also, there is nothing indicating that the first parameter of the `\textcolor` command is the colour parameter and the second is the text to be coloured: neither a human nor an automatic tool, say a preprocessor, can distinguish them without proper knowledge of the `color`  $\text{\LaTeX}$  package. Finally, the use of control sequences and delimiters breaks up text into small chunks causing various problems at later processing stages.

A possible solution is to use texteme properties for colour as well as for escaping. For example, with a `cat-code=12` property the user could mark the backslash as textual content. Of course, more user-friendly property name and values could also be used. The same point can be made for various types of spaces (non-breakable, thin, etc.): instead of using active characters like `~` or control sequences like `\,`, these information can be encoded as orthogonal properties of the same space character. This solution is also far simpler and more intuitive than using Unicode's various *non-breakable space*, *thin space*, *non-breakable* and *thin space*, etc., characters.

**LINGUISTIC TOOLS** One of the advantages of textemes is that the set of available properties is open which, together with modularity, makes it possible to integrate  $\Omega_2$  with new text processing applications. For numerous linguistic problems that bear some relation to typography, such an approach can be fruitful: automatically finding word boundaries in Thai or



Chinese text or distinguishing dots (for abbreviations) and periods (at ends of sentences) in English typography are all complicated linguistic problems that transcend the limits of typesetting engines. If appropriate, standalone linguistic tools already exist and are able to perform the tasks in question, they can be called as external modules of  $\Omega_2$  in order to add linguistic information to the input text in the form of texteme properties. Then, an  $\Omega_2$  developer just needs to implement a second, much simpler module that interprets such linguistic properties and takes them into account at the typesetting stage (correct line breaks for Thai, changes in the widths of spaces for English).

## CONCLUSIONS

The second version of the  $\Omega_2$  system has two new aspects: *texteme support* and *modularity*. Although still in a prototype stage, the basic framework for running external modules has been implemented in  $\Omega_2$  and the underlying text model has also been adapted to the texteme concept. As an addition, the new  $\Omega_2$  outputs a DVI format that contains both characters and glyphs, and with a patched dvipdfmx utility these information can be incorporated into PDF documents that as a result will be able to provide the reader both with formatted output and with the *original* character string. On the input side, a texteme-compliant text editor tool has been developed, allowing users to enter texteme properties into input documents. OpenType support is partially available:  $\Omega_2$  now reads metrics from OpenType files directly. Work is underway for modular GSUB and GPOS support. The authors are also working on moving T<sub>E</sub>X's hyphenation algorithm into an external module, resulting in easier implementation of improved hyphenation tools.

Work that still needs to be done includes full capabilities of texteme input and output: texteme-based

auxiliary (.toc etc.) files and input format, as well as development of various multilingual modules including support for OpenType layout tables. An especially important feature that still needs further development is generation of PDF documents with both character and glyph information added. The already mentioned dvipdfmx tool is a likely candidate for such an extension. The authors kindly welcome contribution from the T<sub>E</sub>X community in these areas.

## REFERENCES

- [1] Jonathan Kew: *The Multilingual Lion: T<sub>E</sub>X Learns to Speak Unicode*. 27th International Unicode Conference.
- [2] Nishikimi Mikiko, Handa Kenichi, Takahashi Naoto, Tomura Satoru: *The m17n Library—A General Purpose Multilingual Library for Unix/Linux Applications*. Asian Symposium on Natural Language Processing to Overcome Language Barriers (2004). <http://www.m17n.org>
- [3] Yannis Haralambous: *New Hyphenation Techniques in Omega 2*. Proceedings of the EuroT<sub>E</sub>X 2006 (Debrecen, Hungary) Conference.
- [4] Yannis Haralambous, Gábor Bella: *Injecting Information into Atomic Units of Text*. ACM Symposium on Document Engineering 2005.
- [5] Yannis Haralambous, Gábor Bella: *Omega Becomes a Sign Processor*. Proceedings of the EuroT<sub>E</sub>X 2005 (Pont-à-Mousson, France) Conference.
- [6] *The FreeType Project*.  
<http://www.freetype.org>
- [7] *International Components for Unicode*.  
<http://icu.sourceforge.net/>



# Making better PDF

---

HARTMUT HENKEL

von Hoerner & Sulger GmbH

Schlossplatz 8, D-68723 Schwetzingen, Germany

hartmut\_henkel@gmx.de

## ABSTRACT

*Under normal circumstances, once the T<sub>E</sub>X part of pdfT<sub>E</sub>X is happy with the input, and all required stuff like fonts and graphics is available, pdfT<sub>E</sub>X produces a valid PDF file. This talk is about PDF output, looking under the hood, showing PDF internals. First the PDF format and structure is shortly described. As pdfT<sub>E</sub>X now supports object streams, this improved compression scheme will be presented as well. Then follows a look into the code that makes up a typical page. What are the operators? How are fonts selected? How is precision of letter placement kept? Here it's interesting to compare e.g. how efficiently different PDF generators put text on a page. In the past, pdfT<sub>E</sub>X has been optimized for compactness of output, and work is ongoing. There are also primitives (e.g. `\pdfliteral`) to put stuff onto a page manually. This will be described, and hints for user's own experimenting will be given. Finally, tools for PDF inspection will be presented.*



# METAPOST developments

TACO HOEKWATER

Elvenkind, Dordrecht

taco@elvenkind.com

## ABSTRACT

*The new release of METAPOST includes some new features as well as a number of bugfixes. The new functionality includes: the possibility to use a template for the naming of output files; support for cmyk and greyscale color models; per-object PostScript specials; the option to generate Encapsulated PostScript files adhering to Adobe's Document Structuring Conventions; the ability to embed re-encoded and/or subsetted fonts; and support for the GNU implementation of troff (groff).*

## INTRODUCTION

Version 0.901 of Metapost was released at BachoT<sub>E</sub>X 2005. It was mostly a bugfix release, that featured and updated manual and the new `mpversion` primitive on top of a set of bugfixes.

At that time, a new version was promised for the autumn. In hindsight, that was overly optimistic. It is now already the summer of 2006, and the feature set for version 1.0 is now finally frozen. It will be released in time for T<sub>E</sub>XLive 2006 and (hopefully) MikT<sub>E</sub>X 2.5), so to an average user not much time will be lost by the delay.

## BUGFIXES

**STABILITY ISSUES** In previous versions of Metapost, the size of the memory array was not stored in mem file. But in Web2c||base systems, the memory sizes are dynamic and the size that should be used by the executable can change depending on the command||line invocation. This discrepancy resulted in a number of painful and unexpected bugs.

- Disappearing specials from the output
- Incorrect error messages
- Unexplained crashes

This problem will be tackled by storing the required minimum memory sizes in the memory dump file. If an unsolvable mismatch occurs, an error message will be issued.

**TURNINGNUMBER** The current (0.9) Metapost executable has a very simple algorithm to calculate the `turningnumber` operation. It simply connects the

path's points using straight segments, adds up all the angles between those segments, and then divides the result by 360. This only works well if the path segments are well-behaved i.e. they do not self-intersect.

This is already an improvement over the old code in the sense that when it is wrong, it is predictably wrong. But it was a temporary measure, and the next version contains completely new code that calculates true curvature for the path segments.

The new algorithm is based on a mailing list discussion between members of the group. It will be slower, but (finally) 100% accurate.

## NEW FEATURES

**FILE-NAME TEMPLATES** The first of the new feature is support for output file-name templates. These templates use `printf`-style escape sequences and are re-evaluated before each shipout. Numeric fields can be left-padded to a user-supplied width by prepended zeroes.

The new primitive is `filenametemplate`, and it is a string-valued command. The syntax is simply:

```
filenametemplate "%j-%3c.eps";
beginfig(1);
  draw p;
endfig;
```

If the file is saved as `test.mp`, then this will create the output file `test-001.eps` instead of `test.1` of previous versions.

A small set of escape sequences are possible, see Table 1 below for details.

To ensure compatibility with older files, the default value of `filenametemplate` is `%j.%c`. If you assign an empty string, it will revert to that default.

Table 1: Escape sequences for `filenametemplate`

<code>%%</code>	A percent sign
<code>%j</code>	The current jobname
<code>%{0-9}c</code>	The charcode value
<code>%{0-9}y</code>	The current year
<code>%{0-9}m</code>	The numeric month
<code>%{0-9}d</code>	The day of the month
<code>%{0-9}H</code>	The hour
<code>%{0-9}M</code>	The minute

**CMYK COLOR MODEL** Support will be added for the industry-standard CMYK color model. In the simplest form this looks like:

```
beginfig(1);
  draw fullcircle withcmkcolor (1,0,0,0);
endfig;
```

To make more flexible use possible, a new type of expression is introduced. A `cmkcolor` is a quartet of numerics that behaves just like the already existing type `color`.

```
beginfig(1);
  cmkcolor cyan;
  cyan := (1,0,0,0);
  draw fullcircle withcmkcolor cyan;
endfig;
```

The new `cyanpart`, `magentapart`, `yellowpart` and `blackpart` allow access to various components of a `cmkcolor` or the CMYK component of an image object.

**GREYSCALE COLOR MODEL** There are only two new primitives for greyscale support: `withgreyscale` and `greypart`. That is because greyscale values are simple numerics.

```
beginfig(1);
  faded := 0.5;
  draw fullcircle withgreyscale faded;
endfig;
```

An image object cannot have more than one color model, the last `withcolor`, `withcmkcolor` or `withgreyscale` specification sets the color model for any particular object.

**RGB COLOR MODEL** Two new aliases for the already existing RGB color model will be added to plain `.mp`. You are requested to use these new keywords `rgbcolor` and `withrgbcolor` when referring to the old color model.

**OBJECT SPECIALS** The new Metapost will support two specials that can be attached to drawing objects. They are output on their own lines, immediately before and after the object they are attached to.

The new drawing options are `withantescript` and `withpostscript`, their arguments simple strings that are output as-is. It is up to the macro writer to make sure that the generated PostScript code is correct.

```
beginfig(1);
  draw fullcircle
    withantescript "gsave"
    withpostscript "grestore";
endfig;
```

**STANDALONE EPS** If `prologues` is set to the value 2, Metapost will generate a proper Encapsulated PostScript level 2 image that does not depend on dvips postprocessing. In this output mode, fonts not be downloaded, but their definition will be handled correctly (see the next paragraph).

Thanks to a detailed set of comments by Michail Vidiassov, this output mode will adhere to Adobe's Document Structuring Conventions. A private PostScript dictionary will be created to reduce the output size for large images.

**FONT RE-ENCODING** If `prologues` is set larger than 1, any used fonts are automatically re-encoded. Their encoding vectors will be included in the output if that needed.

This code is based on the font library used by dvips and pdfTEX. Following in the footsteps of pdfTEX, there are two new associated primitives: `fontmap` file and `fontmapline`. The string-value argument has the same optional flag that is used by pdfTEX:

`none` replace the current font list completely

`+` extend the font list, keep duplicates

`=` extend the font list, replacing duplicates

`-` remove all matching fonts from the font list

```
prologues := 2;
fontmapfile "+ec-public-lm.map";
beginfig(1);
  draw "Helló, világ" infont "ec-lmr10";
endfig;
```

## FONT INCLUSION

Font inclusion is triggered by prologues being equal to 3. Whether or not actual inclusion / subsetting takes place is controlled by the map files. These can be specified using the syntax explained in the previous paragraph.

**GNU GROFF SUPPORT** Version 1.0 of Metapost will have native support for GNU groff, thanks to a set of patches by Werner Lemberg and Michail Vidiassov.

## FUTURE PLANS

The next release after this one is likely to contain the following:

- A option to build metapost as embeddable library instead of an executable.
- The possibility to store drawing objects
- 64-bit internal calculations instead of the current 32 bits.
- Alternative output formats for easier parsing by script backends
- 12-part transform expressions to make it easier for macro packages to implement three-dimensional points.

## WHERE TO FIND METAPOST

*WWW Homepage and portal*

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

*User mailing list*

<http://www.tug.org/mailman/listinfo/metapost>

*Development & sources*

<https://foundry.supelec.fr/projects/metapost>





# Opening up the type\*

TACO HOEKWATER

Elvenkind, Dordrecht

taco@elvenkind.com

## ABSTRACT

*In the near future, pdfT<sub>E</sub>X will gain the ability to use OpenType fonts. This paper explains in broad terms what will be done and why it will be done in that way.*

## INTRODUCTION

There are many reasons why OpenType support is a good thing for pdfT<sub>E</sub>X to have. For one, many of the latest commercial fonts can only be purchased in OpenType format. For another, most of the new OpenType fonts are of higher quality than their older brethren. Also – this is a minor point, but not completely unimportant – everybody else does it. For as long as pdfT<sub>E</sub>X does not support OpenType, it has no chance of being perceived as a viable competitor to those other systems.

Finally and perhaps most importantly, OpenType allows us to escape from TFM format and its many limitations, without the need to invent another special font format that can only be understood by T<sub>E</sub>X and not much else.

## WISHES AND CONSTRAINTS

Of course we want to fit in the new OpenType support in such a way that it behaves in line with the rest of T<sub>E</sub>X:

- It should be possible to use all of the glyphs and features in the font.
- The implementation, its input, and its output should all be platform independent: same font, same syntax, same typesetting.
- We want to make the implementation extensible, just in case someone comes up with OpenType++ in the next decade.
- It should still be possible to define virtual fonts, and even better would be if adjustments to the

font could be made on the fly.

- The interface should be usable by somebody who is not trained as a programmer.
- A small pdf footprint has to be retained, so some sort of font subsetting has to take place.
- Backward compatibility with traditional T<sub>E</sub>X and Postscript fonts has to remain.
- We to provide full control even at the most basic glyph level, not just limited to turning OpenType features on or off

## UNICODE

When dealing with OpenType fonts, adding support for Unicode is more or less implied. Therefore pdfT<sub>E</sub>X had to be extended to handle Unicode input as well as output.

FILE I/O Partial Unicode support already in the current luaT<sub>E</sub>X code base:

- UTF-8 encoded text input and output
- Characters and tokens use 21 bits for storing character information
- Hyphenation patterns are loaded in UTF-8
- The pool file (strings) and buffer are Unicode enabled.

CHARACTERS AND HYPHENATION Now T<sub>E</sub>X does not have a concept of a character: a ‘charnode’ is actually a glyph in a font, together with the font it should be taken from.

To fix this unpleasant intermingling of glyphs with characters, we will extend pdfT<sub>E</sub>X with two new node types:

\*This work is made possible by a grant from the Colorado State University and is part of a project to create a high-end Arabic typesetting engine based on a merge of pdfT<sub>E</sub>X, Aleph and luaT<sub>E</sub>X. The result will be open source and can be merged into future versions of pdfT<sub>E</sub>X.

- A `\font` node is the result of a font selection command by the user.
- A `\unichar` node is the result of tokens being read in.

Not until after hyphenation has taken place will these two types be merged into the traditional `\TeX` ‘char node’.

The old node list and paragraph building routines intertwined ligature building, hyphenation and line breaking.

The new code separates hyphenation from the line breaking decisions, and can be represented in pseudo-code as follows:

```
sub build_nodes {
  while (<tokens>) {
    if (<letter>) {
      store_node(<character>)
    } else if (<font>) {
      store_node(<font>)
    } else if (makes_node) {
      store_node(<token>)
    } else {
      process(<token>)
    }
  }
}

sub line_break {
  while (<nodes>) {
    if(<character>) {
      grow_word(<character>)
    }
    find_hyphens(<word>)
    replace_by_word(<nodes>)
  }
  characters_to_font_glyphs(<nodes>)
  while (<nodes>) {
    if (<glue> or <penalty>) {
      try_break(unhyphenated)
    } else if (<discretionary>) {
      try_break(hyphenated)
    } else {
      store_width(<nodes>)
    }
  }
}
```

This change makes hyphenation completely independent of the current font encoding.

**LANGUAGES** We believe this is a good opportunity to also tackle another traditional problem in `\TeX`: the `\lccode`, `\uccode` and `\sfcode` tables. These tables

contain information that is conceptually part of the current language, and should not be stored in a font attribute.

We want to increase the importance of `\language` codes and attach much more information to language switches.

## FONT LOADING

In current `pdfTeX`, fonts are internally represented as a large storage heap with few dozen auxiliary tables that store various meta-information and pointers into the heap. All of those are global, and implemented as static objects.

While this is very efficient in terms of speed, it is also very hard to alter a font after it has been loaded, and the unification forces all fonts to offer strictly the same interface.

In the new setup, fonts will be loaded under the direct control of lua code, and they will be presented to the typesetting engine as a single lua table for each loaded font. This table will make the font behave much like an object that can be queried and altered directly by the macro programmer, either from `\TeX` macro code (through `\fontdimen`) or from lua code (through callbacks from the typesetting engine).

The binary font loading routines will be written in compiled C code, perhaps by using a separate library like `libfreetype`.

# The making of a (T<sub>E</sub>X) font

TACO HOEKWATER

Bittext, Dordrecht

info@bitttext.nl

HANS HAGEN

Pragma ADE, Hasselt

## ABSTRACT

We want to introduce a new display font to the T<sub>E</sub>X community. The font is a digitization of a series of Duane Bibby drawings, commissioned by Pragma ADE. The digital version for use with ConT<sub>E</sub>Xt is prepared by Bitttext based on scans provided by Pragma ADE.



Figure 1: The first drawing

## INTRODUCTION

At TUG 2003 in Hawaii, Hans Hagen met with Duane Bibby. Hans was looking for some small images to enliven the ConT<sub>E</sub>Xt manuals. A cutout of a very early sketch can be seen in Figure 1, but it was soon agreed that consecutive drawings were going to be an alphabet.

Nothing much happened after that initial meeting, until the beginning of this year when Hans picked up the thread and got Duane started drawing. The alphabet quickly progressed. Starting in a rather naturalistic style like Duane's 'normal' T<sub>E</sub>X drawings, but later progressing toward a much more cartooned style, as can be seen from the drawings in Figure 2.

For easy of use, it was clear that these drawings should ideally become a computer font. Taco agreed to

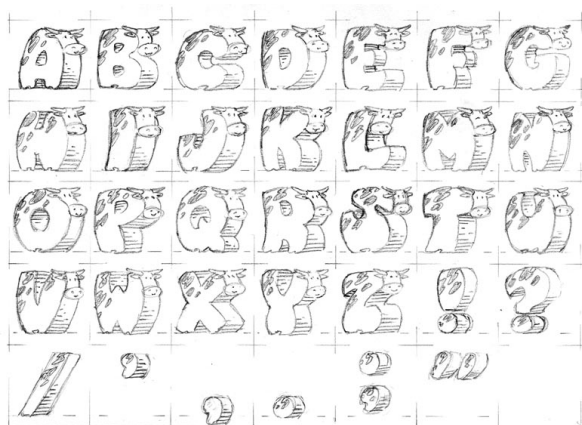


Figure 2: Rough design

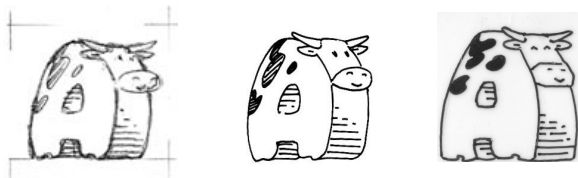


Figure 3: One of the original sheets, showing the alphabet and latin punctuation

take care of the digitization, and luckily the drawings were already prepared for that. As can be seen from the leftmost closeup in Figure 3, the cows are drawn inside a grid. This ensures that they are all the same size, which is a vital requirement for a font.

The center drawing in Figure 3 is a still rather roughly inked version of one of the in-between drawings (there were many). In this particular one you can see that the mouth of the cow was originally more or less oval, but in the final form (on the right) it became much more hexagonal.

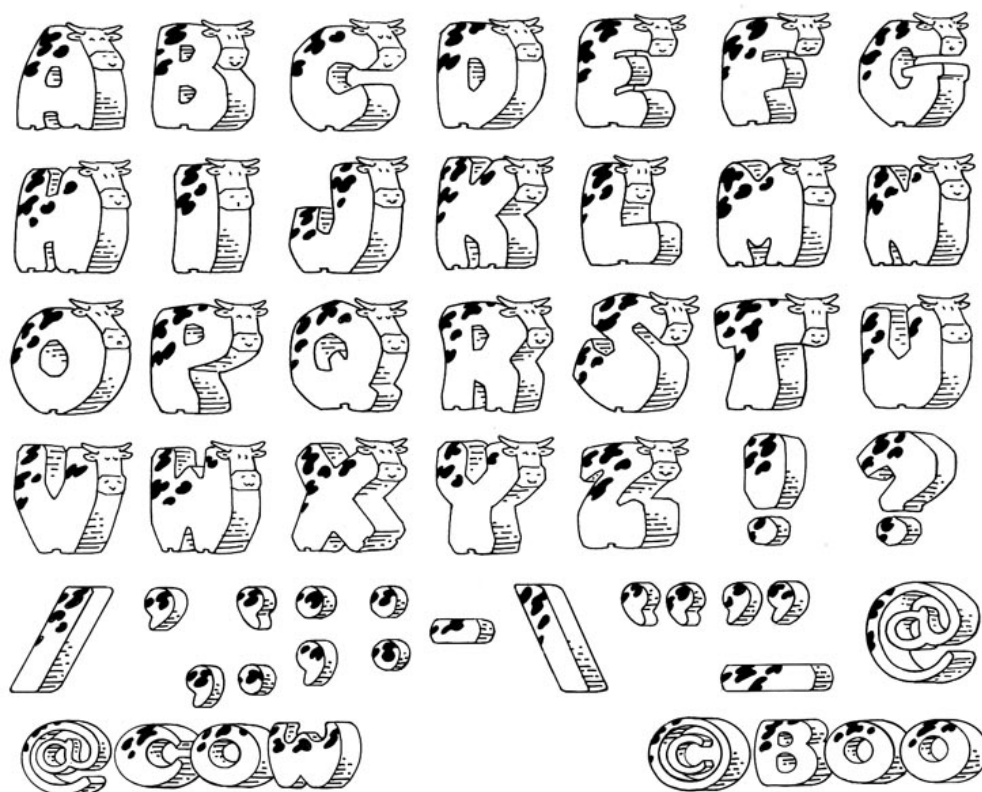


Figure 4: One of the original sheets, showing the alphabet and latin punctuation

## DIGITIZATION

The original sheets were sent to Pragma ADE by regular mail in the beginning of March. Hans scanned the original sheets at 1200 dpi and then forwarded the images to Taco Hoekwater. There were four sheets in all, and one of the is shown in Figures 4. The other three contain T<sub>E</sub>X-related logos and a few (mathematical) symbols.

**PREPARING THE IMAGES** The first task in the preparation of the font was to create a set of bitmap images for use by the import command of fontforge.

For this, the four sheets had to be cut up into many smaller pieces, each containing a single glyph for the to be created font. This being intended as a decorative font, the character set does not even contain the complete ASCII range. Nevertheless, almost a hundred separate images were created.<sup>1</sup>

Fontforge automatically scales imported images so

<sup>1</sup>It would have been nice if this step could have been done solely with free software, but the Gimp turned out to be incapable of handling the 75 Megabyte PNG images for each of the four scanned sheets. Adobe Photoshop was used instead.



Figure 5: One of the imported bitmap images

that they are precisely one em unit high. After cutting the sheets up to pieces, the images therefore had to be adjusted so that they all had the same height. Without that, it would have been near impossible to get all the drawn lines in the glyphs the same width. Figure 5 shows the adjusted version.

**AUTOMATIC TRACING** The autotracer in fontforge (actually the stand-alone autotrace program) does quite a good job of tracing the outlines. But, interestingly enough, only at a fairly low resolution. At

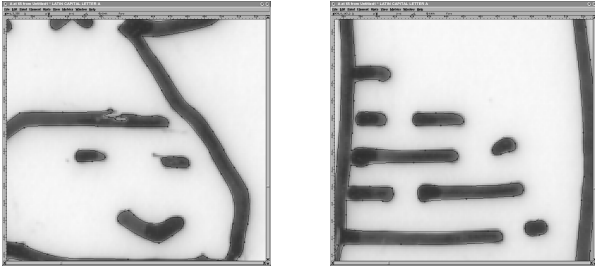


Figure 6: Close-ups of autotracer output

higher resolutions it gets confused and inserts more than a quadratic amount of extra points for each resolution increase. Based on empirical tests, the images were scaled to 40% of their original scanned size, resulting in bitmaps that were precisely 1000 pixels high.

As was to be expected, the autotracer brought out many of the impurities in the original inked version, as you can see in the left image of Figure 6. Luckily, the amount of places where manual corrections like this were needed was not so big as to force us to reconsider the digitization process.

A more severe problem can be seen in the right-hand image of Figure 6. The drawings contain hardly any straight lines. For a font of this complexity, it turned out to be absolutely necessary to simplify the curves. Without simplification, the rendering speed in PDF browsers became unbearably slow. All of the near-horizontal stripes in the bellies were manually removed and replaced by absolute straights.

**HINTING** The final stage in the font editor is to add the Postscript hinting. A screen-shot of a manually hinted letter is visible in Figure 7.

This part of the work is in fact not completely finished yet. It is quite hard to find a balance between two possible extrema.

On the one hand, if there are no hints at all, that results in nice small fonts that render quickly, but ugly.

On the other hand, if there is a lot of hinting information, that creates a much better look but it slows down the rendering. And sometimes it creates even uglier effects, especially with non-commercial renderers.

A middle ground can be reached, but unfortunately only by doing all hinting manually, and that takes quite a lot of time. In the current font, hinting is created automatically by FontLab Studio, except in a few cases where it generated an excessive amount of

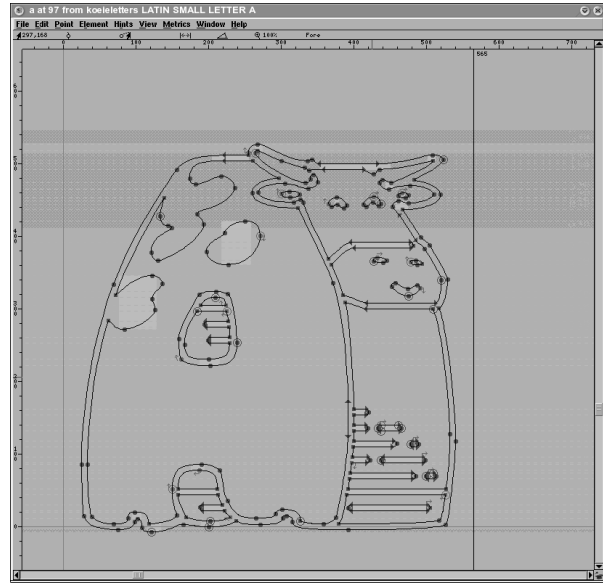


Figure 7: Finished outline

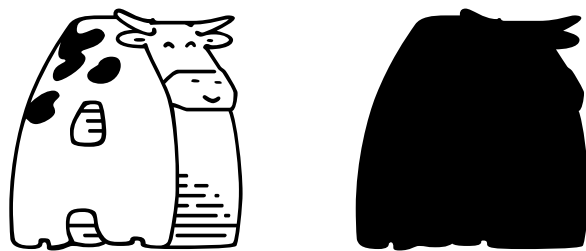


Figure 8: The final 'A'

stem hints.

**FINISHING THE FONT** The font was saved as two separate Postscript Type 1 fonts, one with the text glyphs and one containing the logo glyphs.

The text font is named 'koeieletters', the logo font 'koeielogos'. 'koeieletters' literally translates from dutch to English as 'cowcharacters', but the word 'koeieletter' is also used to indicate a really big character. Like in a billboard, for instance.

Eventually it turned out that we needed a second set of two fonts. Sometimes you want to have text in the cowfont but on top of a colored background. The background would then shine right through the hide of the cow and that was of course unacceptable. Hence, we also have the fonts 'koeieletters-contour' and 'koeielogos-contour'.

Figure 8 shows the final 'A', in the normal and the contour font.



Figure 9: The Pragma ADE logo

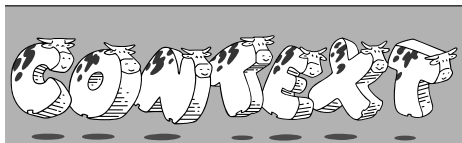


Figure 10: The ConT<sub>F</sub>Xt logo

## PLAYING AROUND

The original goal of this font was to enliven the Pragma ADE manuals. It would be a waste if we did not try to get the most out of the drawings provided by Duane, so we coded some rather silly effects in the font and its metrics. The final paragraphs of this article highlight a few of those.

PRAGMA ADE The only lowercase symbols in the font are in the Pragma ADE logo itself, see on Figure 9.

**THE CONTEXT LOGO** If you look closely at the ConTExT logo on Figure 10, you can see that the shadows and the spots of the cows are drawn in different shades of grey (or colors). This is only possible because there are actually four characters involved instead of just the logo and the background contour, see on Figure 11.

LOGO LIGATURES The line on Figure 12 can be input directly as “Cows in ConTeXt”, thanks to a handcrafted virtual font. This font contains a complete set of ligatures that travel from ‘C at the beginning of a word’, past ‘e following the temporary ligature C\_o\_n\_T’, all the way to ‘t at the end of word’. Such word logos are defined for  $\text{\TeX}$ , ConTeXt and MP, see on Figure 13

MATHEMATICS As is the case for normal text, the math set is also not very extensive, but there are just enough math symbols to allow some example math formulas to be created (see on Figure 14). Virtual fonts make sure that the input is what you expect from  $\text{\TeX}$ .

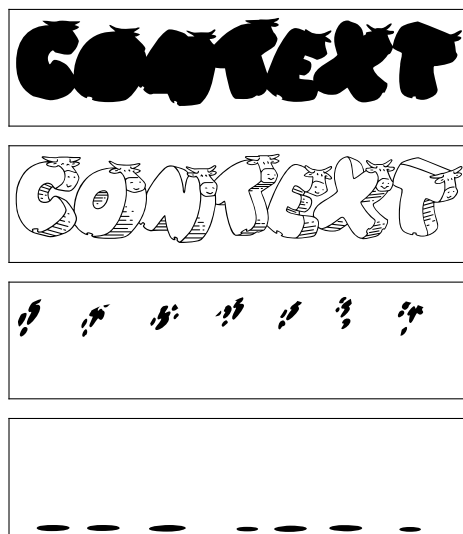


Figure 11: Four characters of the ConT<sub>E</sub>Xt logo



Figure 12: “Cows in ConT<sub>E</sub>Xt” typeset



Figure 13: Word logos

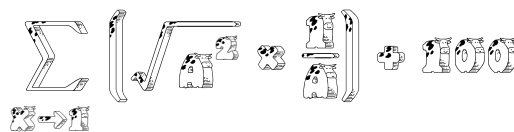


Figure 14: Cows in math mode



Figure 15: Old-style sheep

OLD-STYLE SHEEP Solely seeing cows does tend to get boring after a while. To prevent the font from getting too predictable, we decided we needed some extra freshness. That is why the old-style numerals are actually sheep (see output on Figure 15):

```
I count \oldstylenumerals{10} sheep
```

#### FINAL WORDS

The koeieletter font can be downloaded from the Pragma ADE site. The needed typescript and macros will become part of the standard ConT<sub>E</sub>Xt distribution.





## KöMaL CD – the execution

ILDIKÓ MIKLÓS

KöMaL editorial office

Pázmány Péter sétány 1/A. V. em. 5.106

H-1117 Budapest, Hungary

miklosildiko@t-online.hu

### ABSTRACT

*This talk is about the workflow of production of the CD version of the Hungarian Mathematical and Physical Journal for Secondary Schools using  $\text{\LaTeX}$ .*

### EXTENDED ABSTRACT

The first edition of Mathematical and Physical Journal for Secondary Schools (*Középiskolai Matematikai és Fizikai Lapok* – KöMaL) was appeared on January 1, 1894. From that time several generations of mathematicians and scientists developed their problem-solving skills through KöMaL. The best solutions with the names of the 14–18 year-old authors are printed in the periodical. KöMaL regularly reports on national and international competitions, prints articles on interesting results in mathematics and physics, and includes book reviews. For more than 30 years all the new problems have appeared in English as well as Hungarian in the journal. This means thousands of mathematics and physics problems and exercises in

English. The periodical KöMaL appears in Hungarian language 64 pages a month 9 times a year and in English language two times a year.

From the beginning since 1991 KöMaL was typed by hand. In 1995 we scanned all pages of the issues and made a CD with the pages. We made a searchable database for articles, problems, competitions etc. In 2004 we decided to type all the 36 thousand hand-made pages, and to convert them to MathML. We needed manpower for this project and finally in 2005 we won a tender under European Union. So we could start the job with 24 low educated, handicapped persons. We had only two weeks to teach them to use computer and another 7 days to use  $\text{\LaTeX}$ . In my talk I'll tell this workflow.



# *MLBibTeX meets ConTeXt*

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 reports first experiment of using MLBibTeX – our reimplementation of BibTeX – with ConTeXt, a TeX format more modern than L<sup>A</sup>TeX. We show how to take as much advantage as possible of both ConTeXt and MLBibTeX features when they are used together. Besides, many end-users get used to put L<sup>A</sup>TeX commands inside values of BibTeX fields, and such commands may be unrecognised by ConTeXt. We explain how patterns and preambles allow us to solve such problems.*

**Keywords** ConTeXt, bib module, bibliographies, bibliography styles, BibTeX, MLBibTeX.

## INTRODUCTION

Listing the bibliographical references cited within a document can be done manually – if the L<sup>A</sup>TeX word processor is used, that consists of typing the successive `\bibitem` commands of a `thebibliography` environment [15, § 4.3.2] – but such an approach leads to difficultly maintainable and reusable texts, because they tightly depend on *bibliography styles*. A publisher or anthology's editor would like authors' last names of a 'References' section to be typeset using small capitals whereas another publisher would require the use of standard Roman letters for these last names. Likewise, first names may be abbreviated or put *in extenso*, w.r.t. the bibliography style used... In addition, doing a document's bibliography manually is error-prone: if this bibliography is *unsorted*, that is, if the order of items is the order of first citations of these items throughout the document, some change within the document's body can cause the bibliography to be reorganised. Likewise, keys based on author-date system [16, § 12.3] are to be recomputed if the bibliography is enriched.

A better method consists of using a **bibliography processor**: such a program is given *citation keys*, searches bibliography database files for resources associated with these keys, and arranges them according to a bibliography style, the result being a source file for a 'References' section, suitable for a word proces-

sor. A well-known association between a word and bibliography processor is given by L<sup>A</sup>TeX and BibTeX [18], working in tandem, although this example is not unique. As another example, Tib [1] has sometimes been used with *Plain TeX* [14]; more generally, other examples of bibliography processors are given in [22].

'Historically', BibTeX was initially designed with Scribe<sup>1</sup> to work [21]. In fact, only a few points related to TeX are hard-wired within BibTeX: using braces as delimiters, viewing a group such as `{\command ...}` as an accent command applied to arguments in order to produce a single character [16, § 13.2.2], the use of the `~` character for unbreakable spaces when names are formatted [17, § 5.4], the `width$` function, provided by the style language, that returns the width of a string, expressed using TeX units [16, Table 13.8]. Thus bibliographic entries specified with BibTeX should be usable with any *format* built out of TeX, provided that end-users do not put L<sup>A</sup>TeX-specific command inside field values. Let us recall that TeX basically provides a powerful framework to format texts nicely, but to be fit for use, the definitions of this framework need to be organised in a format. The first formats were *Plain TeX* and L<sup>A</sup>TeX, another format, more modern, is ConTeXt [4]. A biblio-

<sup>1</sup>That is why BibTeX uses the `@` character for specifying its commands and entry types: this character introduces a command name, like `\` in TeX. This convention is still used within Texinfo [2], the GNU documentation format.

```
@BOOK{meaney2003,
  AUTHOR = {John Meaney},
  TITLE = {Context},
  PUBLISHER = {Bantam Books},
  YEAR = 2003,
  NUMBER = 2,
  SERIES = {The \emph{Nulapeiron}
    Sequence},
  NOTE = {The Sequel to ‘Paradox’.
    [Pas de version française
    connue] ! french},
  LANGUAGE = english}
```

Figure 1: Example of a bibliographical entry.

graphic module, usable in conjunction with BibTeX, has been added to ConTeXt [5, 7]. This module defines ConTeXt commands to deal with the components (metadata) of bibliographical information<sup>2</sup>.

These last years, we have designed and implemented a ‘new BibTeX’ – MLBibTeX, for ‘MultiLingual BibTeX’. Of course, it has been designed to work with L<sup>A</sup>T<sub>E</sub>X, but we plan to use it for other output formats than L<sup>A</sup>T<sub>E</sub>X [9]. This article aims to report first experiment of using MLBibTeX to build outputs suitable for ConTeXt<sup>3</sup>. First, we show how ConTeXt can be easily used to pretty-print bibliography database (.bib) files. Then we explain how an interface between ConTeXt and BibTeX can be improved when MLBibTeX is used. This article should be read without any difficulty by any user familiar with L<sup>A</sup>T<sub>E</sub>X and BibTeX: it requires only basic knowledge of ConTeXt and its bib module. It also refers to some notions of XML<sup>4</sup> and Scheme, the implementation language for MLBibTeX<sup>5</sup>, but basically.

## PRETTY-PRINT BIBLIOGRAPHIES

MLBibTeX’s new syntactic features for bibliographical entries are detailed in [8]. Roughly speaking, any .bib file suitable for ‘old’ BibTeX can be processed by MLBibTeX, except that square brackets are syntactic delimiters. Figure 1 gives an example of a bibliographical entry for a book written in English (the value of the LANGUAGE field, handled by MLBibTeX). The value of the NOTE field includes a text to be put down only

<sup>2</sup>If we compare this module to what is provided for L<sup>A</sup>T<sub>E</sub>X, its approach is close to the jurabib package [16, § 12.5.1].

<sup>3</sup>We have used the more recent version of ConTeXt, included in T<sub>E</sub>X Live 2005, available on DVD-ROM.

<sup>4</sup>eXtensible Markup Language. Readers interested in an introductory book to this formalism can consult [20].

<sup>5</sup>The version used is described in [13].

```
<mlbiblio>
...
<book id="meaney2003" language="english">
...
  <note>
    The Sequel to
    <emph emf="no" quotedf="yes">
      Paradox
    </emph>.
    <group language="french">
      Pas de version française connue
    </group>
  </note>
</book>
...
</mlbiblio>
```

Figure 2: XML tree for bibliographical entries.

in French-speaking bibliographies, this text being enclosed by square brackets and labelled by the french language identifier.

As mentioned in [8], the result of parsing a .bib file can be viewed as an XML tree. For example, parsing the meaney2003 entry results in the XML tree sketched in Figure 2. Such a tree can be saved into a file and displayed *verbatim* or handled by tools belonging to XML’s world. ConTeXt provides a way to handle XML texts [19], so it can deal with such files. Figure 3 sketches a pretty-printer for bibliographical entries by means of ConTeXt commands documented in [19, 6]. These bibliographical entries are displayed using MLBibTeX’s syntax. In addition, MLBibTeX’s new syntax for emphasising the parts of person names [18, § 4], based on keywords, is used. For example:

```
AUTHOR = {first => John, last => Meaney},
```

As another pretty-print feature, typographical effects are put into action, whereas keywords, syntactic delimiters and field names are typeset using type writer face. For example, the SERIES field will be rendered as follows:

```
SERIES = {The \emph{Nulapeiron} Sequence},
```

Likewise, an entry is typeset using the typographical conventions of its own language.

If we look at the text given in Figure 3, we notice that the only heavy part concerns language identifiers. ConTeXt uses ISO codes for languages [4, Ch. 7] and can switch to any language by the \language[...] command, without preliminary declaration as L<sup>A</sup>T<sub>E</sub>X needs when the babel package is loaded [16, § 9.2].

```

\enableregime[il1]

\def\ProcessMlBibTeXFieldName[#1]{\tt \expandafter\uppercase{#1} = \textbraceleft}}
\def\CloseMlBibTeXFieldValue{\tt \textbraceright},\par}
\def\ProcessMlBibTeXLanguagePart[#1]{\PutLanguageCommand[#1]{\tt LANGUAGE =} #1,\par}
\def\ProcessMlBibTeXNamePart[#1]{\tt #1 =>}
\def\PutLanguageCommand[#1]{\doifelse{#1}{english}{\language[en]}{%
  \doifelse{#1}{french}{\language[fr]}{\doif{#1}{magyar}{\language[hu]}}}}

\defineXMLenvironment[mlbiblio] \startitemize \stopitemize
\defineXMLenvironment[book] {\item {\tt @BOOK\textbraceleft}\XMLpar{book}{id}{*unkeyed*},%
  \startnarrower[left] \ProcessMlBibTeXLanguagePart[\XMLpar{book}{language}{english}]} {%
  \stopnarrower {\tt \textbraceright}}
...
\defineXMLenvironment[author] {\ProcessMlBibTeXFieldName[author]} \CloseMlBibTeXFieldValue
...
\defineXMLenvironment[first] {\ProcessMlBibTeXNamePart[first]} {, }
\defineXMLenvironment[von] {\ProcessMlBibTeXNamePart[von]} {, }
\defineXMLenvironment[last] {\ProcessMlBibTeXNamePart[last]} {}
\defineXMLenvironment[junior] {, \ProcessMlBibTeXNamePart[junior]} {}

\defineXMLenvironment[asitis] {\tt \textbraceleft} {\tt \textbraceright}}
\defineXMLenvironment[emph] {\XMLifequalelse{emph}{quotedf}{yes}{\tt ‘}{}%
  \XMLifequalelse{emph}{emf}{yes}{\tt \textbackslash emph\textbraceleft}\bgroup\em}{}%
  \XMLifequalelse{emph}{emf}{yes}{\tt \textbraceright}\egroup}{}%
  \XMLifequalelse{emph}{quotedf}{yes}{\tt ’}{}}

\def\GroupMarker{! }
\defineXMLenvironment[foreigngroup] {\tt []} {\tt ] : \XMLpar{foreigngroup}{language}{*error*}} }
\defineXMLenvironment[group] {\startnarrower[left]{\tt []} {%
  {\tt ] \GroupMarker \XMLpar{group}{language}{*error*}} \stopnarrower}
\defineXMLenvironment[nonemptyinformation] {\tt []}\def\GroupMarker{* } {}

\starttext
\processXMLfilegrouped{...}
\stoptext

```

Figure 3: Pretty-printing an XML tree resulting from parsing .bib files.

MLBibT<sub>E</sub>X’s language identifiers are unambiguous prefixes of ways to write in a language, as explained in [11]. For example:

- polish is for the option of the babel package,
- polski is for the polski package [3, App. F],
- pol is for these last two ways, the final choice depends on what a user puts in the document’s preamble,
- po is ambiguous because it may be the prefix of ‘Polish’ or ‘Portuguese’.

In fact, a correspondence table between MLBibT<sub>E</sub>X’s language identifiers and ISO codes, very partially implemented by the \PutLanguageCommand in Figure 3, is needed<sup>6</sup>.

<sup>6</sup>How to put into action our implementation of this table in Scheme

## CONTEXT AND MLBIBT<sub>E</sub>X TOGETHER

If you want BibT<sub>E</sub>X to generate ConT<sub>E</sub>Xt commands for specifying publications from bibliographical entries, you have to use the \setupbibtex command, as explained in [7, § 2.4]. This command gets access to bibliography styles suitable for ConT<sub>E</sub>Xt. Since MLBibT<sub>E</sub>X can process bibliography styles using the bst language [17] in compatibility mode [12], it can deal with these styles. However, we do not recommend this solution, which should be temporary, from our point of view. In addition, the compatibility mode is not very efficient for sake of implementation issues. A first improvement could be the development of new bibliography styles, using the nbst language, close to

is shown in Figure 4.

XSLT<sup>7</sup> [23] and described in [8].

As mentioned in [16, § 13.5.2], the choice among different styles for displaying person names, work titles, ... causes a combinatorial explosion. Besides, all the functions of a bibliography style of BibTeX must be grouped into a unique file, so a rich library of bibliography styles for BibTeX should include a huge number of styles, each being monolithic. As explained in [10], several fragments of a bibliography style written using the nbst language can be assembled dynamically, provided that there is no conflict among rules. Consequently, designing styles according to a modular approach is easier in MLBibTeX than in BibTeX. Moreover, an extended version of the `\setupbibtex` command could allow the use of *several complementary files* for a bibliography style.

## CONTEXTE VS L<sup>A</sup>TEX

In order to increase the expressive power of the information put into .bib files, end-users sometimes put L<sup>A</sup>TEX commands within the values of BibTeX fields. Some basic commands are recognised by ConTeXt, some not. There are two solutions to this problem:

- when outputs for ConTeXt are produced, the contents of `@preamble` rubrics included in .bib files is not put as BibTeX would do; MLBibTeX uses `@contextpreamble` rubrics, so they can be used to implement some L<sup>A</sup>TEX commands in ConTeXt;
- a better solution is given by patterns, expressed using Scheme, replacing substrings by XML-conformant strings; for example<sup>8</sup>:

```
(define-pattern "\\emph{#1}"
  "<emph>#1</emph>")
```

Patterns aim to process any L<sup>A</sup>TEX command put in a .bib file, including user-defined commands, as explained in [9]. But only some predefined patterns are implemented now.

## DIRECT INTERFACE

ConTeXt does not deal with the same auxiliary files than L<sup>A</sup>TEX. Moreover, it builds an .aux file only if the `\setupbibtex` command is activated. Let us recall that BibTeX reads only .aux files, never .tex files.

<sup>7</sup>eXtensible Language Stylesheet Transformations, the language of transformations used for XML texts.

<sup>8</sup>Let us recall that in Scheme, the ‘\’ character is used to escape special characters in constant strings. To include it with in a string, it must be itself escaped.

However, MLBibTeX may need to parse the preamble of a source file, as explained in [11]. Concerning outputs suitable for ConTeXt, the information of interest is the encoding: can MLBibTeX put accented letters of Latin-1 directly<sup>9</sup>, or does it have to use accent commands of T<sub>E</sub>X?

A better solution, making useless .aux files and the parsing of preambles is to build a *driver* written in Scheme. Of course, this task requires some knowledge of both the Scheme programming language and the broad outlines of MLBibTeX’s implementation, but the result is a small-sized program, as shown in Figure 4. You can see how to add a citation key as if was caught from an .aux file; *job-name* is the base name of the file processed by ConTeXt; *filename-list* groups all the fragments of a bibliography style. The `and-let*` macro causes a sequential evaluation to be stopped as soon as a false value (for a failure) is returned. Otherwise, the non-false result may become the value of a local variable (examples are `sxml-biblio-tree` and `k1`).

## CONCLUSION

When we began this task, we had written only some small examples using ConTeXt. And we were afraid of reprogram some important parts of MLBibTeX. To be honest, changes were needed, but not as many as we believed. Concerning the bib module, we learned it quicker than we believed. The first meeting between MLBibTeX and ConTeXt has succeeded.

## ACKNOWLEDGEMENTS

Many thanks to Hans Hagen and Taco Hoekwater, who kindly answered my very ConTeXt-nical questions.

## REFERENCES

- [1] James C. Alexander: *Tib: a T<sub>E</sub>X Bibliographic Preprocessor*. Version 2.2, see CTAN:biblios/tib/tibdoc.tex. 1989.
- [2] Robert J. Chassell and Richard M. Stallman: *Texinfo. The GNU Documentation System. Version 4.8*. <http://www.gnu.org/software/texinfo>. December 2004.
- [3] Antoni Diller: *L<sup>A</sup>TEX wiersz po wierszu*. Wydawnictwo Helio, Gliwice. Polish translation of L<sup>A</sup>TEX

<sup>9</sup>The internal representation uses accented letters as single characters.

```
(and-let* (((log-output-p-pv 'open) job-name)) ; Opening the .log file.
  (((bibtexkey-alist-pv 'add-key) "hoekwater2001")) ; Citation key to be processed.
  (((bibtexkey-alist-pv 'extend))) ; Processing all the entries.
  ((let ((bib-suffix ".bib"))
    (every (lambda (filename) ; Parsing .bib files.
      (s-parse-bib-file (filename-plus filename bib-suffix #f)))
    bibliographyfile-list)))
  (sxml-mlbiblio-tree (s-get-sxml-mlbiblio-tree)) ; Build the XML tree.
  (((language-trie-pv 'use-iso-code-table))) ; Using ISO codes for all the languages.
  (((preamble-pv 'set) "contextpreamble")) ; Using @contextpreamble{...} as preambles.
  (k1 (n-assemble-nstyles filename-list)) ; Styles are assembled and compiled into a function.
  (((output-encoding-pv 'set) 'latin1)) ; Accented letters of Latin-1 allowed in the output file.
  (((bbl-output-p-pv 'open) job-name))) ; Opening the output file.
(k1 sxml-mlbiblio-tree) ; Applying the whole style to the XML tree.
((bbl-output-p-pv 'close)) ; Closing files.
((log-output-p-pv 'close))
#t) ; Final result.
```

Figure 4: MLBib $\TeX$ 's kernel for use with Con $\TeX$ t.

- Line by Line* with an additional annex by Jan Jelowicki. 2001.
- [4] Hans Hagen: *Con $\TeX$ t, the Manual*. November 2001. <http://www.pragma-ade.com/general/manuals/cont-enp.pdf>.
- [5] Taco Hoekwater: "The Bibliographic Module for Con $\TeX$ t". In: *Euro $\TeX$  2001*, pp. 61–73. Kerkrade (the Netherlands). September 2001.
- [6] Taco Hoekwater: *Con $\TeX$ t System Macros. Part 1/ General Macros*. 2002. <http://tex.aanhet.net/context/syst-gen-doc.pdf>.
- [7] Taco Hoekwater: *Con $\TeX$ t. Module Documentation*. March 2006. <http://dl.contextgarden.net/modules/t-bib/doc/context/bib/bibmod-doc.pdf>.
- [8] Jean-Michel Huffle: "MLBib $\TeX$ 's Version 1.3". *TUGboat*, Vol. 24, no. 2, pp. 249–262. July 2003.
- [9] Jean-Michel Huffle: "MLBib $\TeX$ : beyond L $\TeX$ ". In: Karl Berry, Baden Hughes and Steven Peter, eds., *Preprints for the 2004 Annual Meetings*, pp. 77–84. TUG, Xanthi, Greece. August 2004.
- [10] Jean-Michel Huffle: "Making MLBib $\TeX$  Fit for a Particular Language. Example of the Polish Language". *Biuletyn GUST*, Vol. 21, pp. 14–26. 2004.
- [11] Jean-Michel Huffle: *Managing Languages within MLBib $\TeX$* . To appear. June 2005.
- [12] Jean-Michel Huffle: "Bib $\TeX$ , MLBib $\TeX$  and Bibliography Styles". *Biuletyn GUST*, Vol. 23, pp. 76–80. In *Bacho $\TeX$  2006 conference*. April 2006.
- [13] Richard Kelsey, William D. Clinger, Jonathan A. Rees, Harold Abelson, Norman I. Adams iv, David H. Bartley, Gary Brooks, R. Kent Dybvig, Daniel P. Friedman, Robert Halstead, Chris Hanson, Christopher T. Haynes, Eugene Edmund Kohlbecker, Jr, Donald Oxley, Kent M. Pitman, Guillermo J. Rozas, Guy Lewis Steele, Jr, Gerald Jay Sussman and Mitchell Wand: "Revised<sup>5</sup> Report on the Algorithmic Language Scheme". *HOSC*, Vol. 11, no. 1, pp. 7–105. August 1998.
- [14] Donald Ervin Knuth: *Computers & Typesetting. Vol. A: the  $\TeX$ book*. Addison-Wesley Publishing Company, Reading, Massachusetts. 1984.
- [15] Leslie Lamport: *L $\TeX$ . A Document Preparation System. User's Guide and Reference Manual*. Addison-Wesley Publishing Company, Reading, Massachusetts. 1994.
- [16] Frank Mittelbach and Michel Goossens, with Joannes Braams, David Carlisle, Chris A. Rowley, Christine Detig and Joachim Schrod: *The L $\TeX$  Companion*. 2nd edition. Addison-Wesley

Publishing Company, Reading, Massachusetts.  
August 2004.

- [17] Oren Patashnik: *Designing BibTeX Styles*. February 1988. Part of BibTeX's distribution.
- [18] Oren Patashnik: *BibTeXing*. February 1988. Part of BibTeX's distribution.
- [19] PRAGMA ADE, <http://www.pragma-ade.com/general/manuals/example.pdf>: *XML in ConTeXt*. November 2001.
- [20] Erik T. Ray: *Learning XML*. O'Reilly & Associates, Inc. January 2001.
- [21] Brian Keith Reid: *SCRIBE Document Production System User Manual*. Technical Report, Unilogic, Ltd. 1984.
- [22] David Rhead: *The "Operational Requirement" (?) for Support of Bibliographic References by L<sup>A</sup>T<sub>E</sub>X 3*. Technical Report L3-005, L<sup>A</sup>T<sub>E</sub>X 3. August 1993.
- [23] W3C: *XSL Transformations (XSLT). Version 1.0*. W3C Recommendation. Edited by James Clark. November 1999. <http://www.w3.org/TR/1999/REC-xslt-19991116>.



# Appendix G illuminated

BOGUSŁAW JACKOWSKI

BOP s.c., Gdańsk, Poland

B\_Jackowski@gust.org.pl

## INTRODUCTION

This paper aims to provide a collection of illustrations to *Appendix G* of *The T<sub>E</sub>Xbook* [1].

To begin with, I will summarize briefly the main issues of *The T<sub>E</sub>Xbook* which will be dealt with here; next, I confine myself to the explanation of the figures. Obviously, I will use the same notation as is used in *Appendix G*.

I recommend reading this paper simultaneously with *Appendix G*, although they partly overlap.

## MOTIVATION

T<sub>E</sub>X's algorithm for typesetting mathematical formulas is precisely described by Donald E. Knuth in *Appendix G*. The description suffices to implement the algorithm in other languages. For example, it was implemented in JavaScript by Davide P. Cervone [2].

The only drawback of *Appendix G* is that no illustrations are provided. Of course, it is only a relative drawback. Professor Knuth apparently can live without illustrations. My comprehension critically depends on pictures. When they are missing in the original text, I am making sketches while reading.

A few years ago, during an umpteen reading of *Appendix G*, I prepared a bunch of sketches for myself. I didn't think about publishing them, as I was convinced that it is just my predilection or idiosyncrasy; but, judging from the paucity of math fonts around, I concluded that perhaps others might have similar problems.

Therefore, I decided to publish my illustrations to *Appendix G* in hope that they may prove useful, for example, for those working on math extensions for the available non-math fonts. Moreover, they may turn out to be helpful in future works on the improvement of T<sub>E</sub>X; after all, the algorithm is older than a quarter of century, and the world is not sleeping. For example, Murray Sargent III from Microsoft published recently (April, 2006) an interesting note on using Unicode for

coding math [3]. He apparently was inspired by T<sub>E</sub>X: the notation is certainly T<sub>E</sub>X-based and well-known names appear in the acknowledgements and bibliography (Barbara Beeton, Donald E. Knuth, Leslie Lamport).

## MATH STYLES

In math formulas, the following eight styles are used:

- $D, D'$  – in display formulas, generated out of text placed between double dollars  $$$\dots$$$  (*display style*);
- $T, T'$  – in formulas occurring in a paragraph, i.e., placed between single dollars  $\$. \dots \$$  (*text style*);
- $S, S'$  – in formulas occurring in lower or upper indices, i.e., after the symbols  $\wedge$  and  $\substack{\scriptstyle}$  (*script style*);
- $SS, SS'$  – in formulas occurring in indices of indices or deeper (*scriptscript style*).

Typically, in the plain format, 10-point fonts are used for the styles  $D$  and  $T$ , 7-point fonts for the style  $S$ , and 5-point fonts for the style  $SS$ . The “primed” styles, called cramped, use the same point sizes; they differ from the uncramped ones in the placement of subformulas. Table 1 defines the relations between styles of formulas and their subformulas.

The symbol  $C$  will denote the current style, the symbol  $C\uparrow$  – the corresponding style of a superscript, symbol  $C\downarrow$  – the corresponding style of a subscript. From the rules given in Table 1 it follows that  $C\downarrow = (C\uparrow)'$ .

## MATH LIST

T<sub>E</sub>X reads the math material and makes a math list out of it. The math list contains the following math-specific objects:

- *atom*, the basic element;

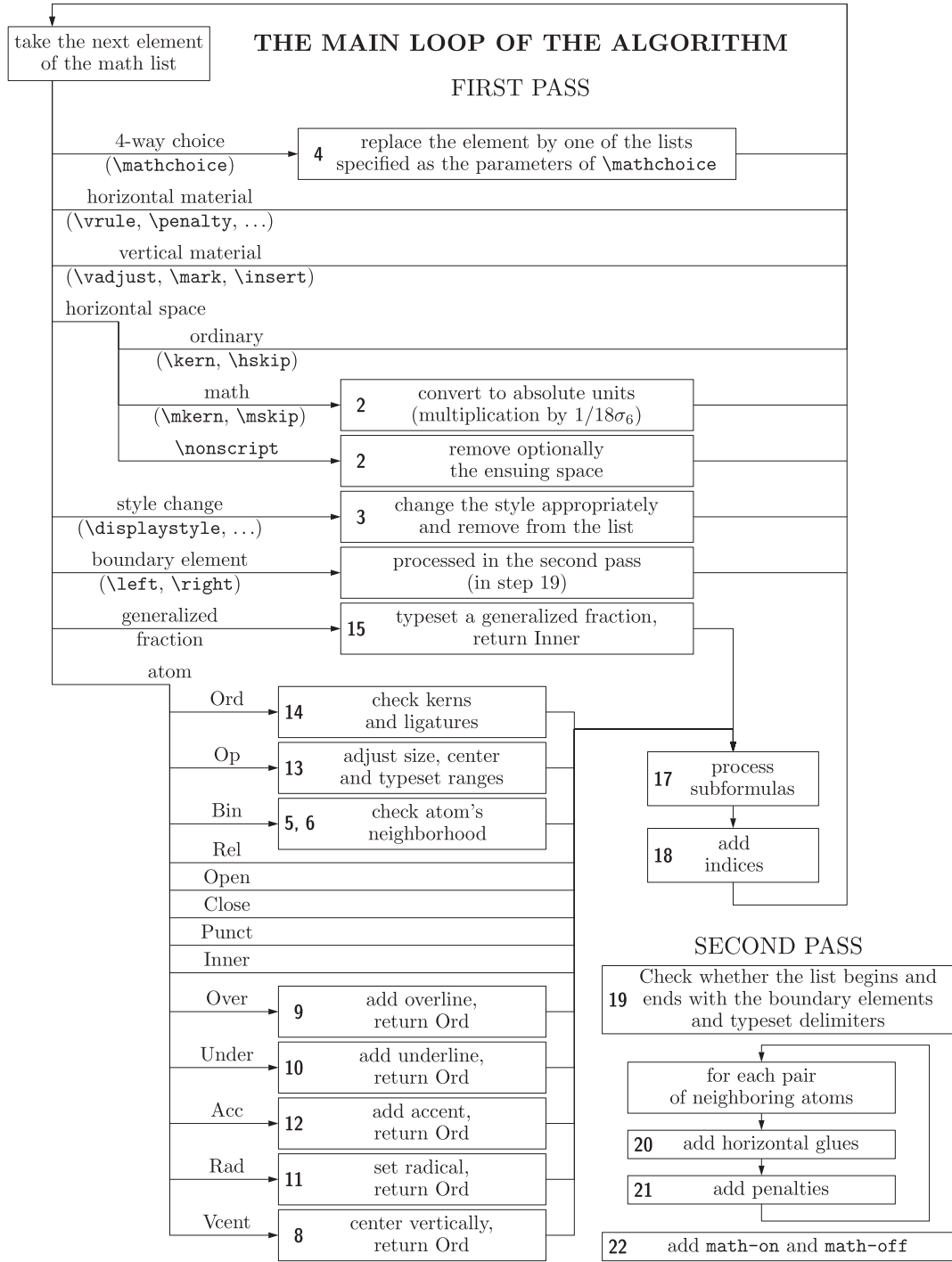


Figure 1: The scheme of the algorithm of the math list processing; the numbers at the left of the boxes refer to the steps of the algorithm, as described in *Appendix G*

Basic style	Superscript style	Subscript style
$D, T$	$S$	$S'$
$D', T'$	$S'$	$S'$
$S, SS$	$SS$	$SS'$
$S', SS'$	$SS'$	$SS'$

Basic style of the formula $\alpha \over \beta$	Numerator ( $\alpha$ ) style	Denominator ( $\beta$ ) style
$D$	$T$	$T'$
$D'$	$T'$	$T'$
$T$	$S$	$S'$
$T'$	$S'$	$S'$
$S, SS$	$SS$	$SS'$
$S', SS'$	$SS'$	$SS'$

Table 1: Rules of the style change – indices (top) and fractions (bottom)

- *generalized fraction*, the result of `\above`, `\over`, etc.;
- *style change*, the result of `\displaystyle`, `\textstyle`, etc.;
- *boundary element*, result of `\left` or `\right`;
- *4-way choice*, the result of `\mathchoice`.

There are several types of atoms, depending on their contents. Table 2 (actually, a duplicate of the table given on the page 158 of *The T<sub>E</sub>Xbook*) lists all possible cases.

Name	Description of the atom
Ord	ordinary, e.g., $x$
Op	holding a big operator, e.g., $\sum$
Bin	holding a binary operator, e.g., $+$
Rel	holding a relational symbol, e.g., $=$
Open	holding an opening symbol, e.g., $($
Close	holding an closing symbol, e.g., $)$
Punct	holding a punctuation symbol, e.g., $,$
Inner	“inner”, e.g., $\frac{1}{2}$
Over	holding an overlined symbol, e.g., $\overline{x}$
Under	holding an underlined symbol, e.g., $\underline{x}$
Acc	holding an accented symbol, e.g., $\hat{x}$
Rad	holding a radical symbol, e.g., $\sqrt{2}$
Vcent	holding a <i>vbox</i> produced by <code>\vcenter</code>

Table 2: The types of atoms which may appear in a math list

Moreover, the math list may contain elements specific for a vertical list:

- *horizontal material* (*rules*, *penalties*, *discretionaries* or *whatsits*);
- *vertical material*, inserted by `\mark`, `\insert` or `\vadjust`;
- *horizontal space*, inserted by `\hspace`, `\kern`, or (acceptable only in math mode) `\mskip`, `\nonscript` or `\mkern`.

The math list is processed twice; after the second pass the “normal” vertical list is created. The scheme of the algorithm (consisting of 22 steps) is shown in Figure 1.

Obviously, font dimension parameters play an essential role. Following *The T<sub>E</sub>Xbook*, I will denote the  $i^{\text{th}}$  parameter of the second family (`\textfont2`, `\scriptfont2`, `\scriptscriptfont2`) by  $\sigma_i$ , and the  $i^{\text{th}}$  parameter of the third family (`\textfont3`, `\scriptfont3`, `\scriptscriptfont3`) by  $\xi_i$ . Note that the  $\sigma_i$  parameters have different values for different styles, while the values of the  $\xi_i$  parameters do not depend on the current style when Computer Modern fonts and the plain format are used.

#### SELECTED STEPS OF THE ALGORITHM

In most cases, the steps of the algorithm are straightforward and a short verbal explanation suffices. Some more detailed elaboration needs, in my opinion, only the typesetting of: radicals (step 11), mathematical accents (step 12), operators with limits (step 13), generalized fractions (step 15), and formulas with indices (step 18). The ensuing subsections deal with these steps.

**TYPESETTING RADICALS** The process of assembling a formula containing a radical is presented in Figure 2. The top part of the picture shows the components (the radicand is typeset using the  $C'$  style), the bottom part – the result of the assembling. Let  $w_x, h_x, d_x$  denote respectively the width, the height and the depth of the radicand, and  $w_y, h_y, d_y$  – those of the radical symbol. The height of the radical symbol is expected to be equal to the thickness of the math rule, i.e.,  $h_y = \theta = \xi_8$ . (A font designer may decide that  $h_y \neq \xi_8$  but I cannot see the rationale for such a decision.) The quantity  $\psi$  is defined as follows:

$$\psi = \begin{cases} \xi_8 + \frac{1}{4}\sigma_5, & \text{styles } D, D', \\ \frac{5}{4}\xi_8, & \text{other styles.} \end{cases}$$

The quantity  $\Delta$  is computed in such a way that the radical symbol is vertically centered with respect to the radicand:  $\Delta = \frac{1}{2}(d_y - (h_x + d_x + \psi))$ .

The baseline of the resulting formula coincides with the baseline of the radicand.

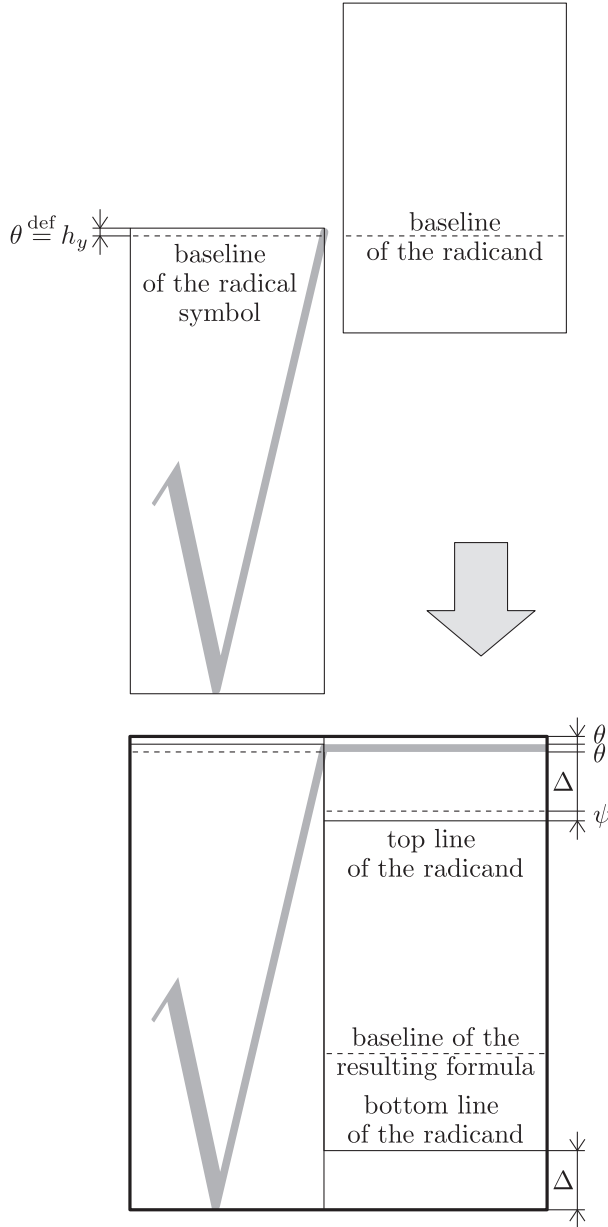


Figure 2: Assembling a radical; symbols are explained in the text

#### TYPESETTING MATHEMATICAL ACCENTS

The typesetting an accented formula is simpler than the typesetting of a radical. Nevertheless, Figures 3 and 4 reveal non-trivial subtleties of the routine.

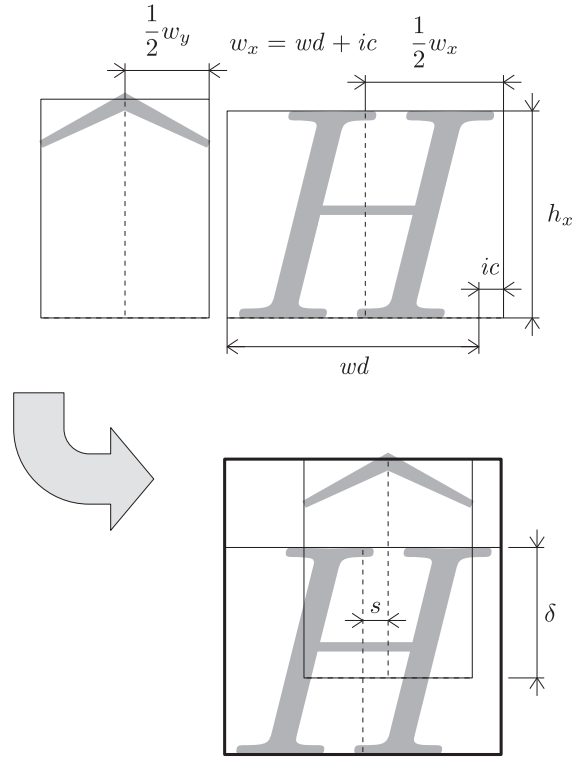


Figure 3: Assembling an accented formula,  $w_y \leq w_x$ ; symbols are explained in the text

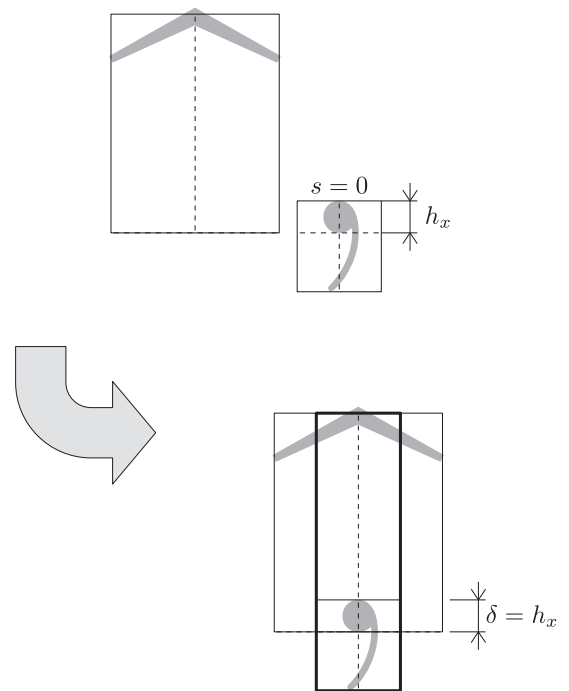


Figure 4: Assembling an accented formula,  $w_y > w_x$ ; symbols have the same meaning as in Figure 3

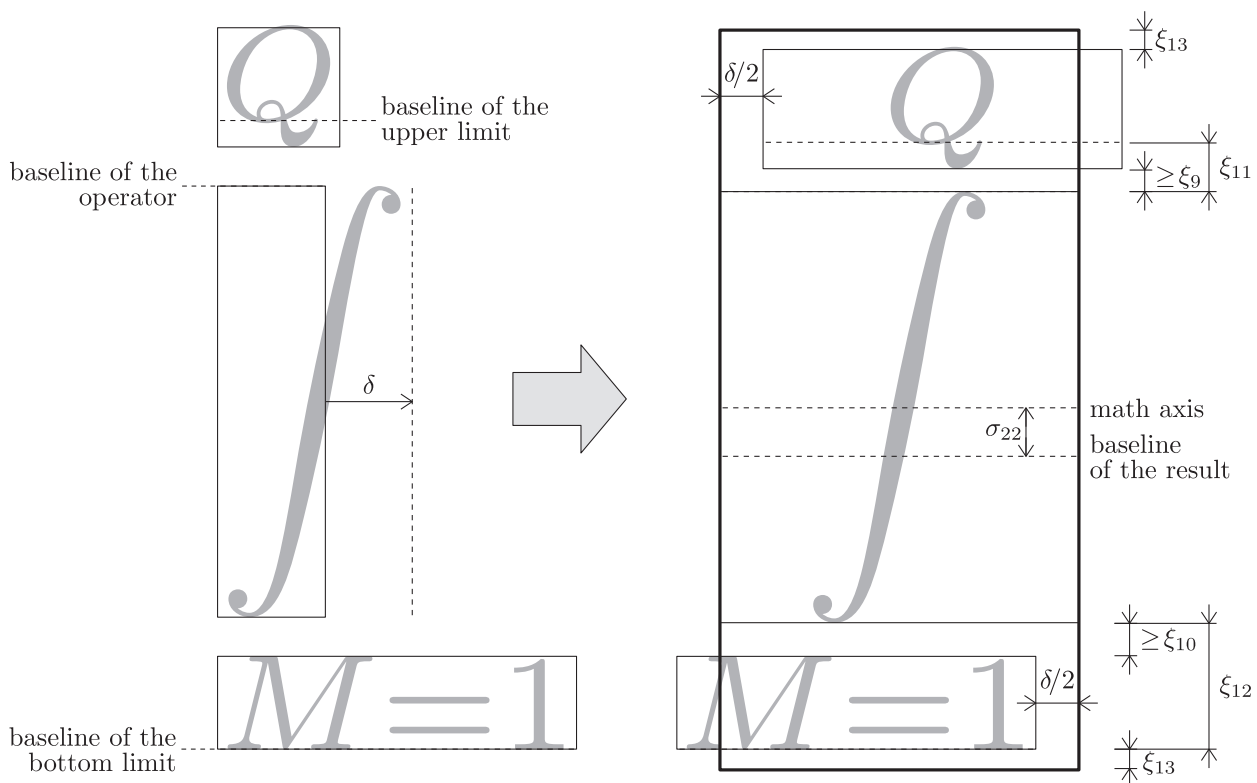


Figure 5: Assembling an operator with limits placed above and below it;  $\delta$  denotes the italic correction of the operator symbol

Like previously, let  $w_x$ ,  $h_x$ ,  $d_x$  denote respectively the width, the height and the depth of the accentee (typeset in the style  $C'$ ), and  $w_y$ ,  $h_y$ ,  $d_y$  – those of the accenter. Actually, these are the widths of the respective boxes; if the accentee is a symbol, its width,  $w_x$ , is computed as the sum  $wd + ic$ , where  $wd$  is the nominal (metric) width of the accentee, and  $ic$  – the italic correction.

Both the accenter and accentee boxes are put into a *vbox* one above the other, and a negative vertical kern,  $-\delta$ , is inserted between the boxes, where  $\delta = \min(x\text{-height}, h_x)$ . The  $x\text{-height}$  is given as the fifth dimen parameter (`\fontdimen5`) of the accenter font.

The horizontal shift of the accenter,  $s$ , is equal to the implicit kern between the accentee and the special character, *skewchar* (defined by the command `\skewchar`); in the plain format, it is the character of code 127 (*tie after*) for family 1, and the character of code 48 (*prime*) for family 2. The kern has nothing to do with the shape of the `\skewchar`, but is intended

to provide an appropriate correction due to the skewness of the accentee. If the accentee is already a boxed formula,  $\text{\TeX}$  assumes that  $s = 0$ .

The width of the resulting formula is always equal to the width of the accentee,  $w_x$ ; the baseline of the resulting formula coincides with the baseline of the accentee.

#### TYPESETTING OPERATORS WITH LIMITS

The placement of the limits of an operator depends on the current style and the usage of `\limits` and `\nolimits` commands. If the style is  $D$  or  $D'$  and the operator `\nolimits` was not applied, the limits are placed above and below the operator, as displayed in Figure 5; otherwise, unless the operator `\limits` was used, the limits are processed as fractions (see section about fractions on the page 83).

The operator symbol is centered vertically with respect to the math axis ( $\sigma_{22}$ ).  $\text{\TeX}$  tries to place the upper formula in such a way, that its baseline is distant by  $\xi_{11}$  from the top of the opera-

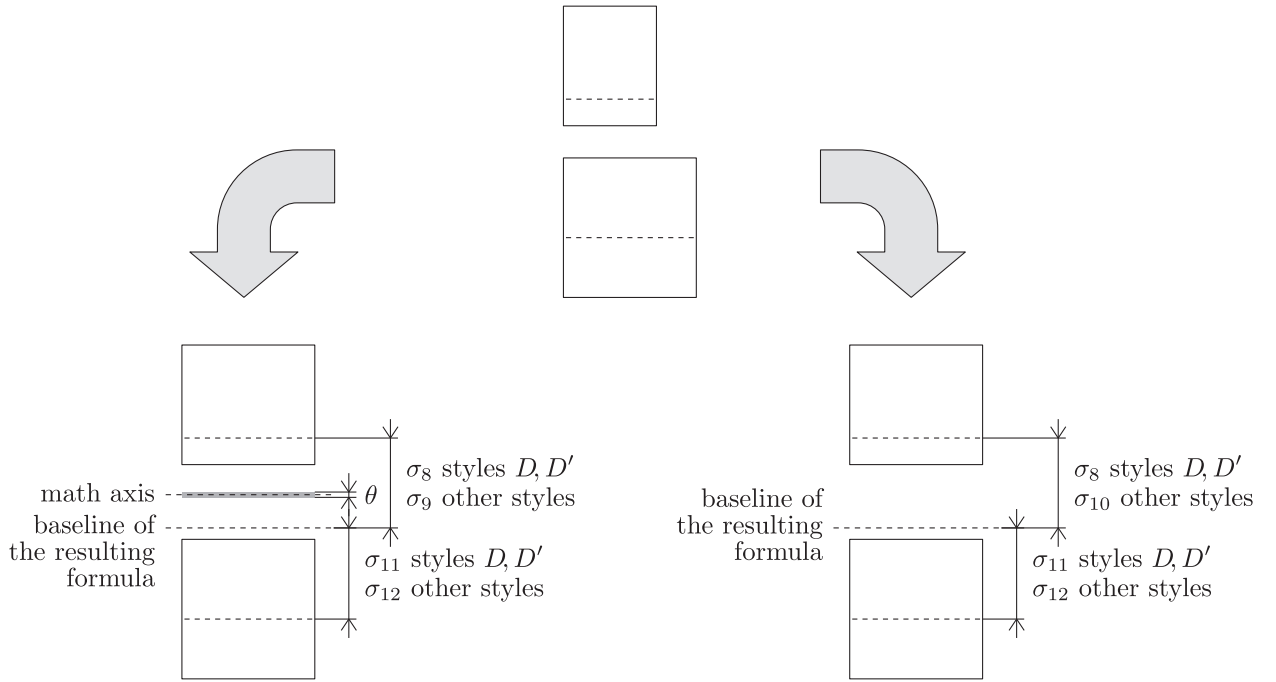


Figure 6: The placement of numerators and denominators in generalized fractions; the thickness of the rule,  $\theta$ , is given either by the value of  $\xi_8$  or explicitly; the latter possibility is provided by the `\abovewithdelims` command; observe that  $\sigma_9$  is used for the formula with a bar, while  $\sigma_{10}$  – for the formula without a bar

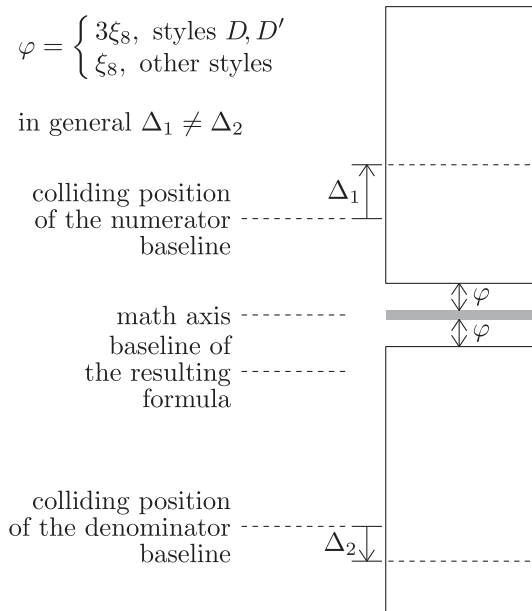


Figure 7: Resolving a collision between the numerator and denominator in the case of a fraction with a bar

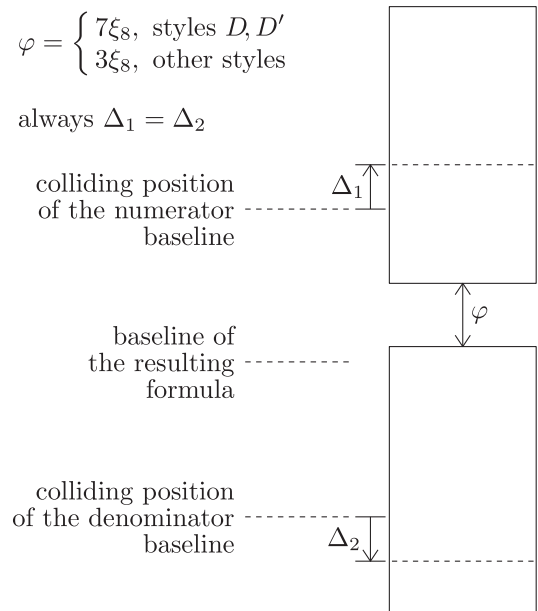


Figure 8: Resolving a collision between the numerator and denominator in the case of a fraction without a bar

tor; however, if the distance between the bottom of the upper subformula and the top of the operator would be less than  $\xi_9$ , the distance  $\xi_9$  is forced. Similarly, the baseline of the lower subformula is distant by  $\xi_{12}$  from the bottom of the operator, unless the distance between the top of the lower subformula and the bottom of the operator would be less than  $\xi_{10}$ , in which case the distance  $\xi_{10}$  is forced.

For the correction of the horizontal placement of the limits, the value of the italic correction of the operator symbol (denoted by  $\delta$  in Figure 5), is used.

### TYPESETTING GENERALIZED FRACTIONS

There are two kinds of fractions implemented in  $\text{\TeX}$ : with or without a bar between the numerator and denominator. They are typeset using different rules, as shows Figure 6. These rules, however, do not suffice, as the numerator and denominator may likely collide.  $\text{\TeX}$  cleverly avoids collisions, as is shown in Figures 7 and 8.

For a fraction with a bar, the numerator and denominator are shifted independently in order to provide a minimal gap,  $\varphi$ , between the formulas and the bar. The position of the bar remains intact – it coincides with the math axis (see Figure 7). For a fraction without a bar, a different strategy is used to avoid the collision, namely, both the numerator and denominator are shifted apart so that the gap between them is equal to  $\varphi$  (see Figure 8). Note that in both cases  $\varphi$  has a different meaning.

### TYPESETTING FORMULAS WITH INDICES

The placement of indices is a fairly complex task due to the variety of situations that may occur.

Figures 9a–9c show the horizontal placement of indices. If the kernel is a single symbol, the superscript, if occurs, is shifted to the right by the amount of the italic correction of the kernel symbol. Technically, a slightly different procedure is involved in the presence of a subscript, as stated in the description of step 17 in *Appendix G*. If the kernel is already a boxed formula,  $\text{\TeX}$  assumes that  $\delta = 0$ . A kern of the value  $\backslashscriptspace$  ( $s$  in the figure) is always appended to index formulas.

The procedure for the vertical placement (Figures 10a–10b) makes use of 7 parameters: from  $\sigma_{13}$  to  $\sigma_{19}$ . Again, different procedures are employed depending

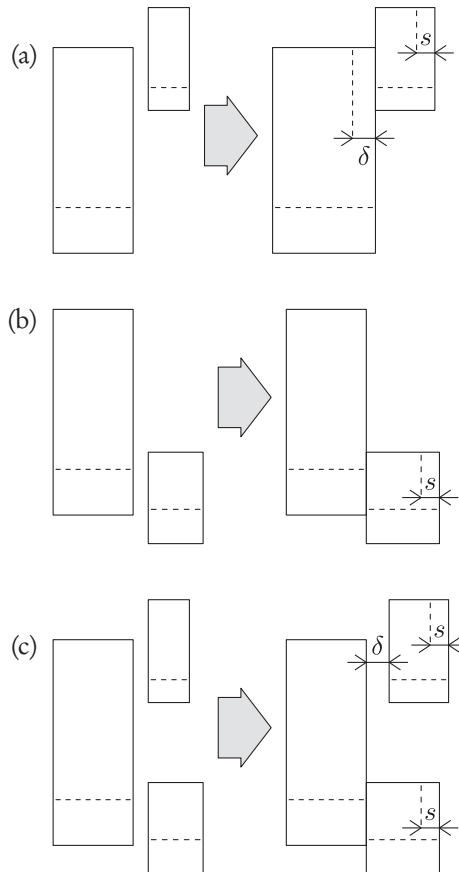


Figure 9: The horizontal placement of indices:

- (a) the placement of a lone superscript,
- (b) the placement of a lone subscript,
- (c) the placement of both superscript and subscript

on the structure of the kernel. If it is a symbol,  $\sigma_{13}$  for the style  $D$ ,  $\sigma_{14}$  for other uncramped styles, and  $\sigma_{15}$  for cramped styles are used for the placement of a superscript; for the placement of a subscript,  $\sigma_{16}$  is used if a superscript is absent and  $\sigma_{17}$  otherwise. If the kernel is a boxed formula,  $\sigma_{18}$  is used for the positioning of the superscript and  $\sigma_{19}$  – for the positioning of the subscript; moreover, the respective values are not taken from the current font:  $\sigma\uparrow$  and  $\sigma\downarrow$  mean that the parameters refer to the fonts corresponding to the styles  $C\uparrow$  and  $C\downarrow$ , respectively.

Indices, as one can expect, are also subject to potential collisions. The actions for such a case are depicted in Figures 11a–11d: (a) the bottom of the superscript formula cannot be placed lower than  $\frac{1}{4}\sigma_5$  above the baseline; (b) the top of the subscript formula cannot be placed higher than  $\frac{4}{5}\sigma_5$  above the baseline; (c) the gap between the bottom

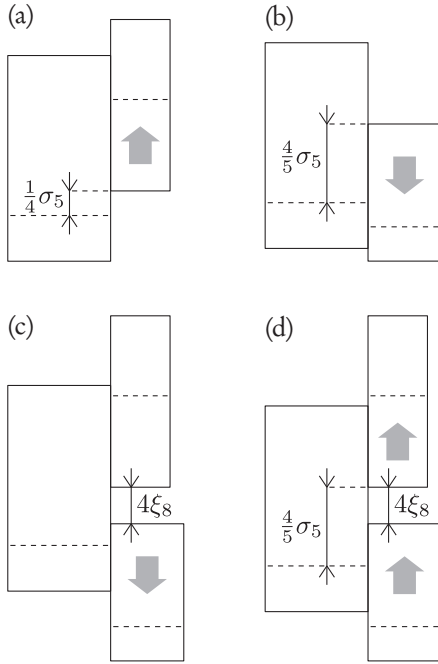


Figure 11: Resolving collisions of indices (further explanations in the text)

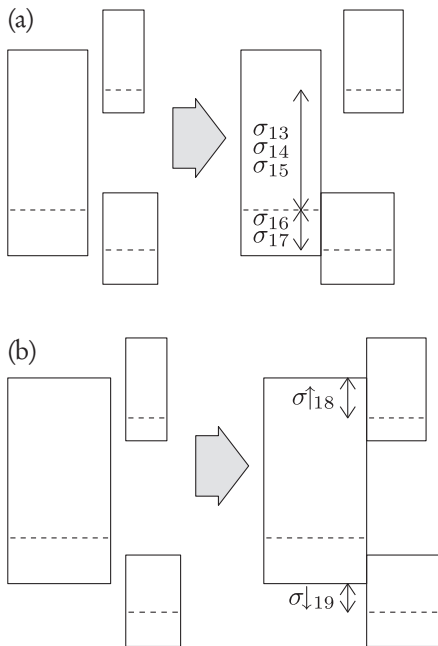


Figure 10: The vertical placement of indices: (a) the placement for the kernel being a symbol, (b) the placement for the kernel being a boxed formula

of the superscript and the top of the subscript cannot be smaller than  $4\xi_8$  (recall that  $\xi_8$  stores the math rule thickness) – the subscript is shifted if required; (d) finally, if the latter situation occurs, both indices can be shifted up so that the bottom of the superscript is not lower than  $\frac{4}{5}\sigma_5$  above the baseline.

## CONCLUSIONS

As was already mentioned, I prepared the illustrations initially for myself and only later decided to publish them in hope that somebody else may also benefit. Therefore, I will eagerly welcome any feedback.

## ACKNOWLEDGEMENTS

I am very grateful to Jerzy Ludwichowski and Piotr Strzelczyk for their prompt and willing help.

## REFERENCES

- [1] Donald E. Knuth, *The T<sub>E</sub>Xbook*, Computers & Typesetting / A, Addison Wesley, 1986
- [2] Davide P. Cervone, *jsMath: A Method of Including Mathematics in Web Pages*, <http://www.math.union.edu/~dpvc/jsMath/>
- [3] Murray Sargent III, *Unicode Nearly Plain-Text Encoding of Mathematics*, <http://www.unicode.org/notes/tn28/>



# Font installation the shallow way

SIEP KROONENBERG

Rijksuniversiteit Groningen Faculteit der Economische Wetenschappen Postbus 800, 9700 AV Groningen, The Netherlands  
siepo@cybercomm.nl

## ABSTRACT

*For one-off projects, you can cut corners with font installation and end up with a more manageable set of files and a cleaner T<sub>E</sub>X installation. This article shows how and why.*

## KEYWORDS

Font installation, afm2pl, afm2tfm, TrueType, pdftex, mapfiles

If you are putting together a flyer or invitation or book cover, then it would be nice if you could without too much trouble test a batch of fonts from your CorelDRAW- or Illustrator cd or your Windows font directory, without polluting your T<sub>E</sub>X installation with a lot of stuff you are never going to use again.

This article takes you through the steps needed to use one or more fonts in one particular document. We won't really install the fonts; we just generate the files that T<sub>E</sub>X needs and leave them where T<sub>E</sub>X will find them, *i.e.* in the working directory. This makes it easy to take the project to another system, and easy to clean things up.

We will primarily use afm2pl to generate .tfm (T<sub>E</sub>X Font Metric) files. Later on, we show the steps required for afm2tfm. Both programs are simpler and much faster to use than the usual choice, fontinst. They create few intermediate or unnecessary files and do their job without virtual fonts. Virtual fonts and fontinst have their place, but sometimes there is no good reason to put up with the inevitable mess.

afm2tfm is available on all major free T<sub>E</sub>X implementations. afm2pl is part of current TeX Live distributions. Note that these programs are needed only to create the necessary font support files for T<sub>E</sub>X; once these files have been created, they can be used on any other system, whether or not it contains afm2pl or afm2tfm.

## AN EXAMPLE

We use a decorative script font Pepita that Adobe bundles or used to bundle with some of its software.

pdftex will need the font outline file epscr\_---.pfb, its T<sub>E</sub>X font metrics file epscr7t.tfm and a map-

file containing an entry relating the two. First, we copy not only epsrc\_---.pfb but also epsrc\_---.afm to the working directory. We need the latter file to generate the .tfm file. Next, we enter the following commands on a command line:

```
afm2pl -p ot1 epsrc_---.afm epscr7t.pl
pltotf epscr7t
```

The extensions .afm and .pl are optional. The first command converts the .afm file to an (almost) human-readable text version of the desired .tfm file. The second command creates the more compact binary version.

Before we can use this font, we must L<sup>A</sup>T<sub>E</sub>X tell about it. We do this with a font family definition file ot1myfontfam.fd:

```
\ProvidesFile{ot1myfontfam.fd}
\DeclareFontFamily{OT1}{myfontfam}{}
\DeclareFontShape{OT1}{myfontfam}{m}{n}{
  <-> epscr7t }{}
```

The prefix ot1 indicates the encoding, which tells which characters occur at what positions. The next section will say more about encodings. The parameters to \DeclareFontShape are successively encoding, family name, weight (*e.g.* bold), shape, font file (without extension) and special options. You can normally leave this last parameter empty. With just one family member, we are not fussy about font characteristics and just pick defaults. We also leave this file in the working directory.

This is the code of our first testfile exabasic.tex, which uses this font:

```
\documentclass{article}
\pagestyle{empty}
\pdfmapfile{=epscr7t.map}
```

```
\newcommand{\fancyfont}%
  {\fontfamily{myfontfam}\selectfont}
```

```
\begin{document}
\fancyfont
Hello, world!
```

```
Accents: \'el\'eve bl\'of \'i;
Kerning: WAV, LTa
\end{document}
```

The `\pdfmapfile` command causes `pdflatex` to read the file `epsr7t.map` which tells `pdftex` how to get the font into the output file. The prepended ‘=’ tells `pdftex` that it should read `epsr7t.map` *in addition to*, not instead of the default mapfile, and that in case of a conflict `epsr7t.map` wins.

Now we are ready to compile `exabasic.tex`:

```
pdflatex exabasic
```

This is the result:

*Hello, world!*  
*Accents: élève blöf i; Kerning: WAV, LTa*

## ENCODINGS

We already made brief mention of encodings. Now is the time to dig a little deeper, because it is a topic that can easily trip you up.

An encoding defines what character corresponds to which number. Only numbers between 0 and 255 are allowed. A `.tfm` file associates character metrics directly with character positions and doesn’t know what position represents what character.  $\TeX$  simply makes assumptions about this correspondence or encoding, and if you disagree with those assumptions then you need to load some macro package or other to tell  $\TeX$  otherwise.

We hope that mainstream  $\TeX$  will eventually move to Unicode, which is a comprehensive encoding of all conceivable characters, including far-eastern alphabets and mathematical symbols. When that happens, we can forget about encodings and also do away with many applications of virtual fonts. There are already some Unicode-based variants of  $\TeX$ .<sup>1</sup>

For a PostScript `.pfb`- or `.pfa` font, character

metrics are stored in a separate `.afm` file. These metrics are associated with characters, not with character positions. Therefore you should specify an encoding to `afm2pl` or `afm2tfm`.<sup>2</sup> The same encoding must also be specified in the mapfile entry. A PostScript font usually has more characters than fit into a single encoding.

A parameter ‘-p `texnansi.enc`’ (where ‘`.enc`’ is optional) means that the encoding should be read from a file `texnansi.enc`. This encoding probably has a different internal name.

**OT1 ENCODING** If you don’t tell  $\TeX$  otherwise, it assumes that you use OT1 encoding. This encoding uses only 128 of the 256 available slots.  $\TeX$  creates missing accented characters from an unaccented base character and a separate accent character. Unfortunately, this interferes with hyphenation. Apart from this, the OT1 encoding has various other oddities, and is best avoided. OT1-encoded fonts often have a  $\TeX$  name ending in `7t`.<sup>3</sup> Note that `ot1.enc` comes with `afm2pl` and is probably not available if you don’t have `afm2pl` on your system.

**T1 ENCODING** T1 is the successor to OT1. It uses all available slots, and has lots of accented characters, also for Eastern European languages. Because the T1 encoding left no room for typographic symbols such as ‘%’ or ‘©’ or ‘f’ you will need to get those from a second encoding of the same font. This second encoding is called TS1 or ‘text companion’.

For most traditional PostScript fonts, some of the accented characters in the T1 encoding aren’t actually present and must be created with virtual font technology from a base character and an accent. Since it doesn’t have to be done by  $\TeX$  itself, this is no obstacle to hyphenation.

Although you can tell `afm2pl` to use T1 encoding, it can’t create composite characters, and such composite characters will be missing unless they are already present in the original font.

T1-encoded fonts often have a  $\TeX$  name ending in `8t`.

<sup>1</sup>Omega and its offshoot Aleph are Unicode-based. Users of Mac OS X may be interested in Xe $\TeX$  ([http://scripts.sil.org/cms/scripts/page.php?site\\_id=nrsi&item\\_id=xetex](http://scripts.sil.org/cms/scripts/page.php?site_id=nrsi&item_id=xetex)), which is built on top of a regular  $\TeX$  installation and lets you use Mac OS X unicode fonts directly with  $\TeX$ .

<sup>2</sup>If you don’t specify an encoding, then you get the encoding from the `.afm` file, which is almost certainly not what you want.

<sup>3</sup>For `afm2pl` and `afm2tfm`, font names have no particular meaning. This is one more difference with `fontinst`. I add encoding postfixes such as `7t` and `8t` to font names just as reminders to myself.

THE TEXNANSI ENCODING *Texnansi* has been introduced by Y&Y, the now-defunct company behind Y&Y $\TeX$ , dviwindow and dvipsone. It combines a good selection of both accented letters and typographic symbols, and normally contains everything you need in a single encoding, at least for Western European languages. *Texnansi*-encoded fonts often have a name ending in 8y.

The package *texnansi* selects the *texnansi* encoding and contains some additional code to smooth out incompatibilities with T1 and OT1.

A TEXNANSI EXAMPLE For this example, we use Augie, a handwriting font from  $\TeX$  Live. These are the commands for generating the .tfm and .map files:

```
afm2pl -p texnansi augie_...afm augie8y.pl
pltotf augie8y
```

This is ly1augie.fd (notice the ly1 prefix):

```
\ProvidesFile{ly1augie.fd}
\DeclareFontFamily{LY1}{augie}{}
\DeclareFontShape{LY1}{augie}{m}{n}{
  <-> augie8y }{}
```

This is the LaTeX code:

```
\documentclass{article}
\usepackage{texnansi}
\pagestyle{empty}
\pdfmapfile{=augie8y.map}
```

```
\newcommand{\fancyfont}%
  {\fontfamily{augie}\selectfont}
```

```
\begin{document}
\fancyfont
Hello, world!
```

```
Accents: \'el\'eve bl\"of \"i;
Symbols:
\textparagraph{} \textdaggerdbl{}
\texttrademark{} \textcopyright
\end{document}
```

and this is the result. Notice the extra symbols. These are absent from the T1 encoding and would have required a text companion font.

Hello, world!  
 Accents: élève blöf î; Symbols: ¶ ‡ ™ ©

## TRUETYPE

Another scalable font format is TrueType, which is supported by pdftex but currently not by dvips. Font metrics are stored in the font file itself. Using True-

Type is somewhat more work; the following commands are required to import a TrueType font such as Trebuchet:

```
ttf2afm trebuc.ttf >trebuc.afm
afm2pl -p texnansi trebuc trebuc8y
pltotf trebuc8y
<edit mapfile to replace .pfb with .ttf>
ttf2afm extracts the metric information from the .ttf file.4
```

afm2pl has no way of knowing that the .afm describes a TrueType font, and guesses that the actual fontfile is trebuc.pfb. Therefore you have to fix the mapfile manually in an editor.

We leave it as an exercise for the reader to write the .fd file and  $\LaTeX$  source for the following example:

Hello, world!  
 Accents: élève blöf î; Kerning: WAV, LTa, WAV, LTa.  
 Symbols: ¶ ‡ ™ ©

## FONT-BASED UPPERCASING AND LETTERSPACING

afm2pl comes with an uppercased version *texnanuc* of *texnansi*. Uppercasing, e.g. in headings, works best in combination with letterspacing. For this, afm2pl has a parameter ‘-m’.

*Warning.* afm2pl implements letterspacing with kerns. Unfortunately, the .tfm format can contain only a limited number of kerns. If there are too many in the .pl file then all kerns and ligatures will be dropped from the generated .tfm file! So use this feature with care. fontinst implements letterspacing by adding sidebearings via virtual fonts, and doesn’t suffer from this limitation.

We can create a letterspaced, uppercased version of Trebuchet with the following commands:

```
ttf2afm trebuc.ttf >trebuc.afm
afm2pl -p texnanuc -m 100 trebuc trebucup8y
pltotf trebucup8y
<edit mapfile to replace .pfb with .ttf>
```

A fontfamily and fontshape declaration might look as follows:

```
\ProvidesFile{ly1trebuc.fd}
\DeclareFontFamily{LY1}{trebuc}{}
\DeclareFontShape{LY1}{trebuc}{m}{upp}{
  <-> trebucup8y }{}
```

The fontshape upp for uppercasing is not an official

<sup>4</sup>This will result in an empty encoding, unless you specify an encoding parameter. But we are going to ignore the encoding in the .afm anyhow.

L<sup>A</sup>T<sub>E</sub>X shape but that doesn't seem to matter. You can use the font as follows:

```
\documentclass{article}
\usepackage{texnansi}
\pagestyle{empty}
\pdfmapfile{=trebucupp8y.map}

\begin{document}
\fontfamily{trebuc}\fontshape{upp}
\selectfont Letterspaced uppercasing
\end{document}
```

and this is the result:

LETTERSPACED UPPERCASING

## A FONT FAMILY

The next example uses a real font family, consisting of the usual four family members plus our letterspaced font. So we will need not only `trebuc.ttf`, as in the previous example, but also `trebucbd.ttf`, `trebucit.ttf`, and `trebucbi.ttf`. For each of these we'll have to run the `ttf2afm - afm2pl - pltotf` sequence, and we'll have to edit each of the generated map files, or create a combined mapfile.

This is its code of the `.fd` file:

```
\ProvidesFile{ly1trebuc.fd}
\DeclareFontFamily{LY1}{trebuc}{}
\DeclareFontShape{LY1}{trebuc}{bx}{n}{
  <-> trebucbd8y }{}
\DeclareFontShape{LY1}{trebuc}{m}{n}{
  <-> trebuc8y }{}
\DeclareFontShape{LY1}{trebuc}{bx}{it}{
  <-> trebucbi8y }{}
\DeclareFontShape{LY1}{trebuc}{m}{it}{
  <-> trebucit8y }{}
\DeclareFontShape{LY1}{trebuc}{m}{upp}{
  <-> trebucupp8y }{}

```

And this is the L<sup>A</sup>T<sub>E</sub>X code using it:

```
\documentclass{article}
\usepackage{texnansi}
\pagestyle{empty}
% better combine these mapfiles!
\pdfmapfile{=trebuc8y.map}
\pdfmapfile{=trebucbd8y.map}
\pdfmapfile{=trebucit8y.map}
\pdfmapfile{=trebucbi8y.map}
\pdfmapfile{=trebucupp8y.map}

\begin{document}
\fontfamily{trebuc}\selectfont
Hello, \textbf{world!}
```

Accents: \el\`eve bl"of \i;

Kerning: WAV, LTa, {\it WAV, \textbf{LTa.}}

Symbols:

```
\textparagraph{} \textdaggerdbl{}
\texttrademark{} \textcopyright
```

```
\fontshape{upp}\selectfont
```

Letterspaced uppercasing

```
\end{document}
```

This is the result:

Hello, world!

Accents: élève blöf ï; Kerning: WAV, LTa, WAV, *LTa*.

Symbols: ¶ ‡ ™ ©

LETTERSPACED UPPERCASING

## USING DVIPS

If you go the dvips route, then you cannot use the `\pdfmapfile` macro. Instead, you have to enter additional mapfiles on the command line:

```
dvips -u +mapfile dvifile
```

The prefix `+` to the mapfile parameter is analogous to the `=` prefix for the `\pdfmapfile` macro: it tells dvips to use the named mapfile *in addition to* the default one.

## USING AFM2TFM

The intention of `afm2tfm` is not to create fonts which are used directly by T<sub>E</sub>X. Instead, they serve as a basis for virtual fonts, *i.e.* recipes to compose fonts from other fonts. But it is not too difficult to subvert this intention:

```
afm2tfm epscr___ -T texnansi \
  -v indirect.vpl direct.tfm >direct.map
#rm direct.tfm
vptovf indirect.vpl
rm indirect.vf
<edit direct.map>
```

Note that the `.afm` filename comes *before* the options.

`vptovf` generates two files from `indirect.vpl`: `indirect.vf` and `indirect.tfm`.

You should remove `indirect.vf`, otherwise the dvi driver or pdftex would think that `indirect` is a virtual font.

Normally, you would also remove `direct.tfm`, but I keep it to show you the difference with the file `indirect.tfm`.

Mapfile information is written to standard output, which therefore had to be redirected, as shown above. It contains the following string:

direct PepitaMT

" TeXnANSIEncoding ReEncodeFont " <texnansi  
(everything on one line). This has to be changed into:  
indirect PepitaMT

" TeXnANSIEncoding ReEncodeFont "  
<texnansi.enc <epsr\_\_\_.pfb  
(type it into one line).

The example below displays differences in spacing between the two. *Note.* This is not an example for copying.

*Direct*

*Accents: élève blöf i; Kerning: WAV, LTA*

*Indirect*

*Accents: élève blöf i; Kerning: WAV, LTA*

#### OTHER OPTIONS OF AFM2PL AND AFM2TFM

With both programs you can artificially slant, narrow and widen a font. afm2tfm can also generate artificial smallcaps. Such manipulated fonts rarely look good, though.

afm2pl also has some options for manipulating the ligkern table and for setting spacing parameters. For casual use, you don't bother with these.

#### OPENTYPE

We are seeing more and more OpenType fonts, which are Unicode-based. These consist of either PostScript/

Type 1 or TrueType outlines inside a TrueType wrapper. OpenType fonts may contain huge charactersets, sometimes including smallcaps and oldstyle figures.

OpenType fonts with Type 1 outlines, which have .otf extension, can be converted with `otftotfm`, part of Eddie Kohler's LCDF Typetools and included in T<sub>E</sub>X Live.

OpenType fonts with TrueType outlines have an extension .ttf and can be treated just like TrueType fonts.

#### AD HOC OR GENERIC SOLUTIONS?

Various people have written scripts to automate font installation. ConTeXt users will be familiar with `texfont`, which, by the way, has an option to use `afm2pl` instead of `afm2tfm`.

Each example took several commands on a command line. So why not a script?

I don't install fonts all that often. I like to decide case by case how to do it: what tools to use, what variants to generate, where to install or not install, how to name the fonts...

Under these circumstances, the simplest and best solution is to either do it by hand, or to write little ad-hoc scripts and keep them with the project.



# Managing a network T<sub>E</sub>X installation under Windows

SIEP KROONENBERG

Rijksuniversiteit Groningen Faculteit der Economische Wetenschappen Postbus 800, 9700 AV Groningen, The Netherlands  
siepo@cybercomm.nl

## KEYWORDS

MikT<sub>E</sub>X, TeXnicCenter, filename database, registry, graphic file formats

This paper is about the MikT<sub>E</sub>X installation I maintain for the Economics Department of the Rijksuniversiteit Groningen. We have long been the home of 4T<sub>E</sub>X. But when development of this project stopped, the time came that this T<sub>E</sub>X installation had to be replaced by something else. This something else was going to be MikT<sub>E</sub>X with TeXnicCenter as editor and front end.

There are various Windows editors which support MikT<sub>E</sub>X, *i.e.* editors which have menu items and buttons for compiling and viewing your TeX documents with MikT<sub>E</sub>X. Configuring *e.g.* TeXnicCenter or Winedt for MikT<sub>E</sub>X is almost automatic. TeXnicCenter is free, both in the beer- and in the speech sense. The MikT<sub>E</sub>X site lists a few more free editors. LaTeXEditor<sup>1</sup> and Texmaker<sup>2</sup> seem to have a similar focus as TeXnicCenter.

MikT<sub>E</sub>X itself comes with a configuration program, MikT<sub>E</sub>X Options or `mo.exe`<sup>3</sup>, which has to be started from outside the editor. Over time, the MikT<sub>E</sub>X installation has been accumulating some additions, especially for handling graphics; see further on.

## MOVING FROM 4T<sub>E</sub>X TO MIKT<sub>E</sub>X

I didn't try to create a unified 4T<sub>E</sub>X-style interface for everything, and also didn't try to replicate the functionality of 4T<sub>E</sub>X, but I did collect the local macros, fonts and graphics from 4T<sub>E</sub>X which were still in use and put them into the MikT<sub>E</sub>X installation, sometimes with some minor tweaks.

I dropped support for the old L<sup>A</sup>T<sub>E</sub>X 2.09 since it would have meant real work for something that might not even get used.

The MikT<sub>E</sub>X installation was put online early in

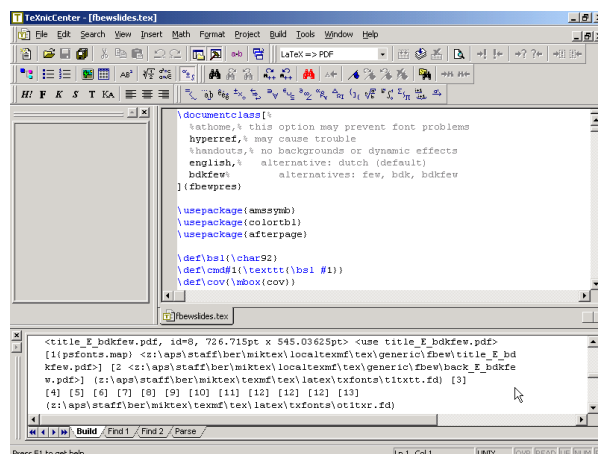


Figure 1: TeXnicCenter, a MikT<sub>E</sub>X frontend

2003. For a year and a half, MikT<sub>E</sub>X and 4T<sub>E</sub>X were available side by side, but in the end, after six months notice, I removed 4T<sub>E</sub>X from the network.

## LAYOUT AND CONTENTS OF THE INSTALLATION

**TEXMF TREES** As to the organization of macros, fonts and other support files, MikT<sub>E</sub>X is rather similar to other modern free T<sub>E</sub>X implementations: it organizes its files into *texmf trees*, which have a standardized structure: *e.g.* font-related files are in subdirectories `fonts\tfm`, `fonts\type1` etc., and L<sup>A</sup>T<sub>E</sub>X macros are in `tex\latex`. Each tree follows such a structure and has its own filename database. Users can configure in which order the trees are going to be searched.

I configured the following three trees: a main tree for files coming from the MikT<sub>E</sub>X distribution; a department tree for local additions, including the files inherited from the 4T<sub>E</sub>X installation; and a user tree for people's own macros and other files. Users have write access only to their own user tree. With this setup, a

<sup>1</sup><http://www.ntu.edu.sg/home5/pg03053527/latexeditor/>

<sup>2</sup><http://www.xmlmath.net/texmaker/>

<sup>3</sup>MikT<sub>E</sub>X also has a package manager. But of course that is not useful to users who don't have write access to the installation.

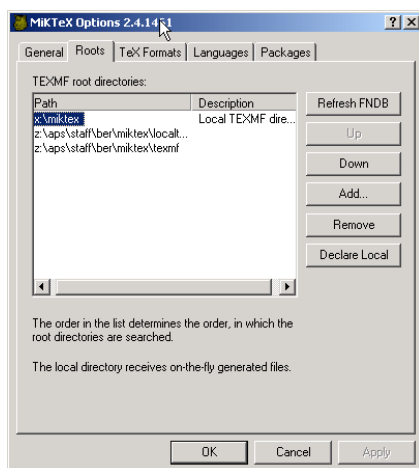


Figure 2: Defining trees and their priorities with MikTeX Options

MikTeX upgrade won't interfere with the department tree, and anything done to the network installation leaves user files alone.

**PACKAGE SELECTION** MikTeX is distributed as a setup program and a set of packages. The MikTeX Setup Wizard lets you choose between three initial package sets: small, large or everything, which you can modify later on. I picked the 'large' package set, and added and removed some packages afterwards.

**ADD-ONS** Non-MikTeX add-ons included originally Ghostscript and GSview, but for the current edition, this was no longer necessary, since Ghostscript and GSview were already installed separately.

For better scripting, I included the Perl .exe- and .dll files, but placed them outside the search path. If people have a need for Perl then they can install their own copy, without these two files getting in the way.

## GRAPHICS SUPPORT

**DRAW PROGRAMS** Our MikTeX installation includes a couple of draw programs. One of these is Ipe (<http://ipe.compgeom.org>), which has a few interesting features:

- It uses pdf<sub>l</sub>atex in the background for typesetting text elements. You can tune typographic details with L<sup>A</sup>T<sub>E</sub>X preamble commands.
- It can import arbitrary pdf via a separate conversion utility.
- A drawing can be layered in the sense that it can be

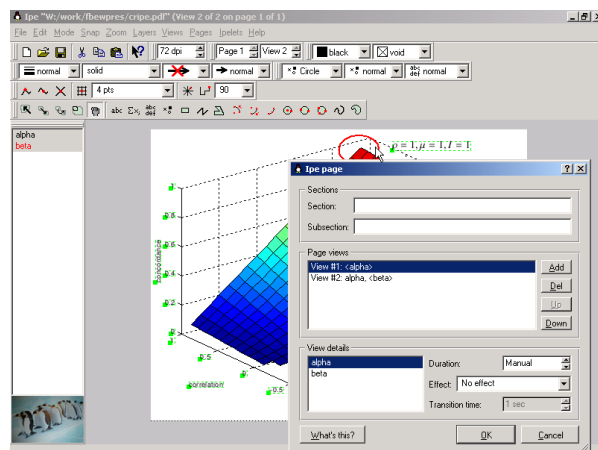


Figure 3: The Ipe draw program has views or pages which are really assemblies of 'layers'.

displayed incrementally in a pdf presentation. In fact, Ipe also advertises itself as a tool for making presentations.

A second draw program is L<sup>A</sup>T<sub>E</sub>XCAD, which generates L<sup>A</sup>T<sub>E</sub>X picture environments. It is very old and basic and is only included for people who still have old L<sup>A</sup>T<sub>E</sub>XCAD drawings in use.

**CONVERTERS** I also added some PostScript-, .eps and .pdf conversion scripts, with desktop shortcuts which can be used as drag-and-drop targets. For conversion from pdf to PostScript I added the xpdf utilities.

I plan to write a basic GUI tool, custom tailored for our installation, which offers all available conversions from a single interface. Of course, the real work will be done by Ghostscript and other trusty command-line standbys.

There is also an installer for wmf2eps. This program offers a fairly practical way to make graphics from MS Office and other Windows programs available to L<sup>A</sup>T<sub>E</sub>X. It seems that not much has happened with it lately, but it still works well enough. It relies on a virtual printerdriver, and therefore isn't a good candidate for a network install.

Its main advantage over simply printing to an eps file is that it calculates a tight boundingbox, rather than just turning an entire page into a graphic.

Recent versions of Ghostscript have a 'bounding-box output device'. There are now various scripts which use Ghostscript to fix boundingboxes. The as



yet nonexistent GUI tool mentioned above should also offer such an option, as an alternative to using `wmf2eps`<sup>4</sup>.

## AN INSTALLER

Installation of the network version involves:

- storing configuration information in the registry
- creating the user tree if it doesn't exist
- creating desktop- and start menu shortcuts
- generating the filename database

Because our desktop systems are very standardized, none of this requires user input.

**FILENAME DATABASE** MikT<sub>E</sub>X's `texmf` trees have an important and annoying difference with the more unixy varieties: the filename database of a tree is not stored with that tree, but in a designated 'local' tree which receives generated files. This local tree can only be the user tree. That means that it is up to the user to update his filename database if items are added to or moved around in the global installation.

The filename database can be generated from the MikT<sub>E</sub>X Options program but, fortunately for writers of installation scripts, it can also be done from the command line:

```
initexmf --update-fndb
```

It happens often enough that the installation fails because the user becomes impatient with the generation of the filename database. Without the help of this database, MikT<sub>E</sub>X becomes very very slow. To minimize the problem, I saved filename database generation till last in the installation process and preceded it with dire warnings about not interrupting the installer. If these warnings are ignored, then the filename database can still be generated after the fact from the MikT<sub>E</sub>X Options program. An alternative would be to copy a pregenerated filename database to the right place during installation.

## DEALING WITH THE REGISTRY

Under Windows, almost any configuration information is stored in the registry.

The registry contains a set of hierarchically orga-

nized keys. There are several root keys. The most important ones are `HKEY_CURRENT_USER` and `HKEY_LOCAL_MACHINE`, `HKCU` and `HKLM` in short. The `HKCU` part of the registry may be on a network drive. `HKLM` is normally in a subdirectory of the Windows directory.

The actual information is contained in *values* under those keys. Values are name-data pairs.

For users, there are very few reasons to edit the registry directly. There are almost always specialized menu entries and dialogs available, such as Tools menus and Control Panel entries.

**REGISTRY TOOLS** If you do need to view or manipulate the registry directly, Windows has a number of tools: you can browse the registry with `regedit`<sup>5</sup>. You can export and import registry keys to and from `.reg` files with either `regedit` or the command-line tool `reg.exe`. These `.reg` files are editable text files. `Reg.exe` may not be installed by default, though. Type '`reg /?`' for help.

Various programming languages, including Perl, VBScript and installer programs, also have functions for accessing the registry.

**FINDING OUT WHAT REGISTRY VALUES ARE NEEDED** There are a number of techniques for identifying the registry settings that you need: one way would be to inspect the source code of the original installer.

A second method is browsing and searching the registry for likely strings, and then testing whether you have captured enough to make the program work as intended. However, such testing can be time-consuming since you sometimes have to re-login or reboot before the changes take effect.

A third method is to export the registry to a text file before installation and after installation or first use, and compare the differences. There exist automated installers which do just this, but the GNU diff program works just fine.

You still have to decide what are the differences that matter. There will probably a lot of spurious differences. For example, most programs record windows positions and their most recently used files in the registry.

<sup>4</sup>There is also a Linux program called `wmf2eps`, but I have had mixed results with it. It seems better to convert `wmf`- and `emf` graphics to `.eps` or `.pdf` on the original system *before* trying to use them elsewhere.

<sup>5</sup>Under earlier versions of Windows, `regedit` and `regedt32` each could do things that the other couldn't. Under Windows XP, `regedit` combines the functionality of the earlier `regedit` and `regedt32`. Its version of `regedt32` simply starts up `regedit`.

Also, some information which occurs only once can appear to occur multiple times. In particular, under Windows 2000 and later, the keys under HKCR are copies from keys under HKLM\software\classes and HKCU\software\classes.

#### ALL USERS OR NOT: HKLM AND HKCU

Often, when you install software, there is a choice whether or not to install for all users. If you do, keys are added under HKLM\software; otherwise under HKCU\software.

**WHERE TO LOOK** The most important settings are under keys `software\<program name>` and `software\classes` (either from HKLM or from HKCU). The keys under `classes` define file types and define what happens when you double-click a file in Windows Explorer. Command-line programs may not need any entries here.

Uninstall information can be found under HKLM\software\Microsoft\Windows\CurrentVersion\Uninstall. MikTeX doesn't have an uninstaller yet.

#### REGISTRY ENTRIES FOR MIKTEX ITSELF

MikTeX makes very modest use of the registry. It just records the locations of the texmf trees, stores uninstall information, and defines the .dvi filetype, associating it with the yap previewer.

I also added the MikTeX binaries directory to the search path, for those people who prefer to run MikTeX from the command line. On Windows 2000 and Windows XP the search path and other environment variables are stored in the registry; for the current user under HKCU\software\environment, for the local system under HKLM\system\currentcontrolset\

control\session manager\environment.

#### REGISTRY ENTRIES FOR GHOSTSCRIPT AND

GSVIEW Ghostscript needs to record the location of the main dll and of its own fonts and support files. GSview defines the .eps- and .ps file types and associates them with itself.

#### REGISTRY ENTRIES FOR TEXNICCENTER

TeXnicCenter stores a lot of information in the registry, but it can configure itself when it is started for the first time if it can find the MikTeX-, Acrobat-, Ghostscript- and GSview registry settings. All the user has to do is to answer 'yes' a few times. I decided to

leave configuration of TeXnicCenter to itself.

It is possible to rerun the TeXnicCenter configuration wizard at a later date. This may come in handy whenever MikTeX or Ghostscript or GSview has moved, or the Adobe Reader has been upgraded.

It would have been nice if TeXnicCenter checked at startup whether these programs are still at their previous location.

#### MORE INSTALLERS

I started out with one installer, but now there are several.

#### A NETWORK INSTALLATION FOR A LATEX

**COURSE** A second network installation was needed for a computer course for econometrics students. This installation is a slightly stripped-down version of the first one: no department tree, and without some of the add-ons.

**A CD INSTALLATION** Earlier, I had already made a cd with on it the standard installers for MikTeX, Ghostscript, GSview and TeXnicCenter, and a copy of the department tree, plus a file with instructions how to put everything together.

There were two problems with this: it was complicated enough that some people preferred to let me install MikTeX for them, and other people figured that they might as well download and install MikTeX directly from the internet. Which was not exactly wrong, but differences in their setup sometimes made it difficult to debug their problems.

So I hope that the new installer cd has persuaded some people not to do a do-it-yourself install.

**THE DIFFERENCES** I already listed some of the differences between staff- and student installations. Some differences between network- and cd installers are:

- With the cd installer, users can choose locations for the main installation and for their own data. These locations are fixed in the network version.
- The cd installer has to copy everything to the harddisk, whereas in the network version everything is already in place. In fact, re-running the network installer is no big deal. Which is just as well, since time and again configurations get messed up by a malfunctioning network or other mishaps.

- The cd installer tests for Ghostscript and GSview. If they are missing, the user first has to install these programs, *e.g.* with the installers provided on the cd. The network version simply knows that Ghostscript and GSview are present and where they are.
- The cd does an ‘All users’ install; the network version doesn’t. Since on our network the HKCU part of the registry and the start menu are on the user’s network drive, you can run MikT<sub>E</sub>X from any workstation on the network.
- The cd only contains MikT<sub>E</sub>X fonts; not the additional department fonts. Adding fonts in MikT<sub>E</sub>X 2.4 is tricky at best. Adding them to systems that I didn’t control caused too much grief.

It was not difficult to create the installer script as one main script with four different wrapper scripts.

I kept the installation and the installer on a Linux Samba server. I managed to put all ‘real’ files in a single directory tree, and to access these files through four different sets of symlinks. This prevented worries about keeping the versions in sync.

**INSTALLER PROGRAMS** The standard way to distribute applications at the university is to create entries in NAL, or Novell Application Launcher, using Novell ZENworks. As I understand, ZEN identifies file system and registry differences before and after installation. With ZEN, an installer can make system changes which the user wouldn’t have permissions for without ZEN. However, a first attempt to create such a NAL entry for MikT<sub>E</sub>X, done together with our NAL specialist, wasn’t exactly plain sailing.

I needed something that I could develop and test on my own system. This was even more important for the student install for the L<sup>A</sup>T<sub>E</sub>X course, where I had to do everything through intermediaries who weren’t even in the same building.

In the first edition, which didn’t include a cd counterpart, I could make do with a batchfile with some embedded Perl code<sup>6</sup> for manipulating the search path.

The cd version of the second edition required user interaction, first for telling users to install Ghostscript if it wasn’t found, second for asking users where MikT<sub>E</sub>X should be installed. So it was really time to switch to a GUI installer.

<sup>6</sup>That is, the batchfile calls Perl with the `-x` switch and itself as parameter; see the perlrun manpage.

I picked NSIS<sup>7</sup>. It is completely scriptable and can be used from the command line<sup>8</sup>. It has functions for reading and writing the registry and for creating shortcuts. It offers string handling and conditionals. You can choose to what extent you want to package files into the installer itself, *i.e.* you can also tell the installer to copy files straight from the installation media to the target system.

The principal drawback of NSIS is very low-level string handling, which is quite painful if you are used to Perl string handling and regular expressions.

I have also heard good things about InnoSetup (<http://www.jrsoftware.org/isinfo.php>), but by then I was almost finished with my NSIS installer.

## DEVELOPMENT AND TESTING

**VIRTUAL MACHINES** Nowadays, you don’t need physical test machines anymore; with software such as VMware you can create virtual guest machines for testing inside a host, *i.e.* inside your everyday PC. The hard disk of this guest computer is a very large file on the host’s disk. Its screen can take up the entire physical screen, but it may also run inside a window, whatever is convenient. If host and guest have similar processors, performance can be quite decent.

VMware has versions for Windows- and for Linux hosts. There are other options for virtual machines, both commercial and free: the Xen project (<http://www.xensource.com>) is getting a lot of attention, and may in time become a very interesting alternative. See also QEMU ([www.qemu.org](http://www.qemu.org)), Win4Lin (<http://win4lin.com>) and Virtual PC (<http://www.microsoft.com>)<sup>9</sup>. Some of these emulators are focused more on running Windows applications than on creating a realistic test environment.

For testing installers, you can create a ‘clean’ virtual machine with just Windows and some indispensable programs installed. Then you can run simulations on copies of this virtual machine. Getting another clean test machine is just a matter of making a fresh

<sup>7</sup><http://nsis.sourceforge.net/>

<sup>8</sup>For GUI addicts, there is an interface with some buttons to push. There are also third-party editors with a GUI for building dialog boxes.

<sup>9</sup>Virtual PC was bought from Connectix by Microsoft in the second half of 2003. The Macintosh version of Virtual PC was at the time the only real option for running Windows on the Mac. I, along with many other Mac owners, was duly shocked by this sell-out. But in the meantime, other emulators appeared, and now that the Mac is moving to Intel, we can hope for VMware-quality virtual machines on Mac OS X from other companies than Microsoft.

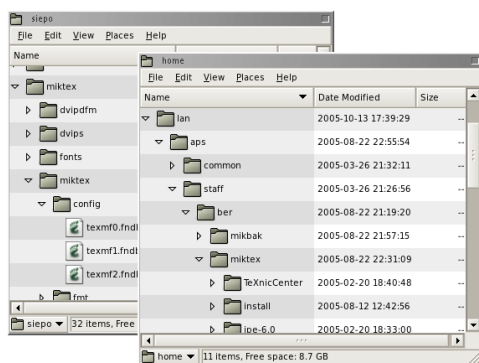


Figure 5: The Samba shares seen from the Linux server. It makes no difference whether the server is a separate machine, a VMware host or a second VMware guest machine.

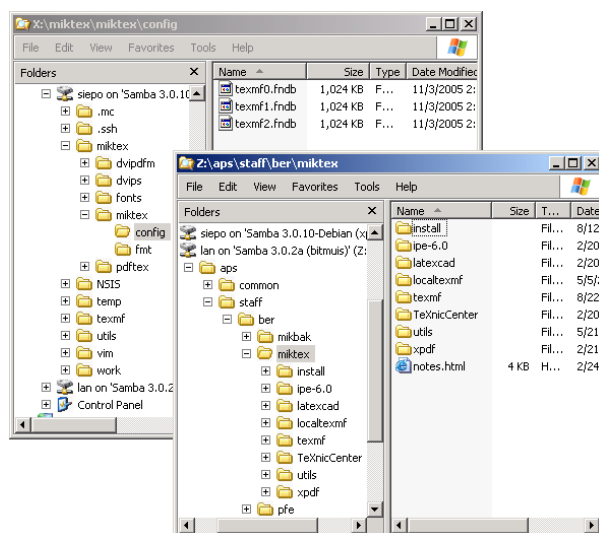


Figure 4: The Samba shares seen from the Windows client machine. It doesn't make a difference whether the Windows machine is physical or virtual.

copy of the original one, which takes only minutes.

**VIRTUAL NETWORKING** For networking, I let VMware create what it calls a host-only network, with no direct access to an external network. This saved me the hassle of protecting virtual Windows machines against malware from the internet. I configured the Linux host as a Samba server, with the MikTeX instal-

lation in a Samba share, and the user's home directory in another share.

**ROAMING PROFILES** The university started using *roaming profiles*. The idea is to place user configuration data as much as possible on their own network home drive. These user configuration data include e.g. their start menu and the HKCU part of the registry.

With Samba, roaming profiles means configuring the (or a) Samba server as a PDC or Primary Domain Controller. This is no fun. It means creating things called machine accounts for the client machines, and painstakingly reading Samba documentation. A very helpful and funny guide was 'The Unofficial Samba HOWTO'. You can find this document via the Samba site. Its current location is <http://hr.uoregon.edu/davidr1/docs/samba.html>.

For testing, I now make clean MikTeX-free profiles, with just the worst default Windows settings fixed, and work with copies of those clean profiles, just as I already did with guest machines.

**DISAPPEARING FILE TYPES** In theory, with MikTeX and TeXnicCenter, roaming profiles should work perfectly: there should be no need to install or configure anything on the machine itself. In practice, definitions of filetypes under HKCU, i.e. all keys and values under HKCU\software\classes, got lost in between logins – in real life, not in my test setup.

For staff network installations, a workaround is to duplicate the filetype definitions in HKLM\software\classes. For student network installations, students don't have write permission on HKLM keys, so this would only be possible with something like ZEN. As a hacky alternative, I put together a registry patch with filetype definitions, i.e. a .reg file with the registry keys and values under the HKCU key which define these filetypes. By double-clicking this file, these keys and values are imported into the registry. Students have to run this registry patch before each MikTeX session.

Next time around, if the disappearing filetypes problem still isn't solved, I might have to do something with ZEN after all.

# Automatic discretionary hyphenation in OpenOffice.org

LÁSZLÓ NÉMETH

Budapest University of Technology and Economics  
Department of Sociology and Communication  
Centre for Media Research and Education  
nemeth@mokk.bme.hu

## ABSTRACT

*This paper presents a new automatic discretionary hyphenation algorithm, which uses discretionary hyphenation patterns, and which is now included into OpenOffice.org.*

## EXTENDED ABSTRACT

A new automatic discretionary hyphenation algorithm has been introduced in OpenOffice.org 2.0.2. Similar to Sojka's work [2] it is an extension of Liang's hyphenation algorithm [1], but it doesn't use special discretionary hyphenation table and characters. Instead, it uses discretionary hyphenation patterns, freely combined with the original non-discretionary hyphenation patterns. An example (only with discretionary hyphenation patterns) in Hungarian language:

```
l.1l/l=1
a1atje./a=t,1,3
e1etje./é=t,1,3
.schif1fahrt/ff=f,5,2
drucker/k=k,4,2
d1dzsel./dzs=dzs,1,4
.as3szon/sz=sz,2,3
n1nyal./ny=ny,1,3
.til1lata./ll=1,3,2
```

And the result in OpenOffice.org [3].

The new algorithm can handle complex discretionary hyphenation: there are automatic Hungarian discretionary hyphenation in OpenOffice 2.0.2. Robust work and simple implementation of the algorithm may be promising also for T<sub>E</sub>X.

## REFERENCES

- [1] Liang, Franklin Mark. *Word Hy-phen-a-tion by Com-put-er*. Stanford University. 1983.
- [2] Sojka, Petr. *Notes on compound word hyphenation in T<sub>E</sub>X*. TUGboat., 16(3), 1995.
- [3] Discretionary hyphenation algorithmt included into OpenOffice.Org.  
<http://www.openoffice.org/nonav/issues/showattachment.cgi/35655/README.discretionary>



# Font Verification and Comparison in Examples

KAREL PÍŠKA

Institute of Physics, Academy of Sciences  
182 21 Prague, Czech Republic  
piska@fzu.cz

## ABSTRACT

*The contribution demonstrates several techniques of verification and comparison fonts widely used with T<sub>E</sub>X: METAFONT fonts and outline fonts in the PostScript Type 1 and OpenType formats. The aim of generating various proofsheets files in PDF or PS with node and control points, control vectors and hinting zones for the subsequent visual scanning of graphic glyph representation, calculations of differences between metric data (e.g. character widths), between contour curves for different versions or releases, etc., is to accomplish the auditing process to be more quick and efficient. Numerous tools — METAFONT, METAPOST, (pdf)(E)T<sub>E</sub>X, dvips, gv, FontForge, MetaType1, T<sub>E</sub>Xtrace, mfttrace, tutils, awk, sed, sort and other programs — are used. The differences “greater than negligible” often indicate problems with compatibility, sometimes they may signal a bug undetected even for a long time. The examples are mostly taken from the current T<sub>E</sub>XLive 2005. The results of verification of CM, EC, LM, CS and other fonts available from T<sub>E</sub>XLive or CTAN, comparison for compatibility and consistency and the information about differences and bugs will be reported.*

## INTRODUCTION

A short version of the article (low level font oriented and technical) is presented here. After a brief explanation of font elements important for typesetting with T<sub>E</sub>X (metrics and glyph images) we will show a limited number of illustrations. The motivation of the work was to prepare tools and intermediate results in a textual (lists, tables) and a visual form to find, detect and demonstrate differences, mistakes, cases of inconsistency and incompatibility and, in the next step, to improve past, current or future fonts.

## FONT TYPES AND FONT DATA

*TFM character dimensions.* The tfm (T<sub>E</sub>X font metric) files contain four dimensions for each character [METAFONT book]

- *charwd*, the width
- *charht*, the height above the baseline
- *chardp*, the depth below the baseline
- *charic*, the character’s “italic correction”

They define the size of each character’s “bounding box” T<sub>E</sub>X needs to typeset. For formatting text T<sub>E</sub>X uses only metric information and does not need glyph shapes. The tfm files also contain an information

about *ligatures* and *kerning* pairs defining the space adjustment between two adjacent characters.

The dvi output produced by T<sub>E</sub>X contains only references to glyphs. The real glyphs are absent in dvi and are included into the final output in PS/PDF by device drivers (e.g. dvips) or by pdfT<sub>E</sub>X. The tfm files can be converted by the tftopl program to human-oriented property list files. We can read, edit and process them in this text form in an easier way. afm is a metric format for Type 1 PostScript fonts. The metric data in tfm and afm are represented and stored differently. The afm bounding box has the different meaning and in afm the glyph width is defined by the WX parameter.

ec-lmr10.tfm/pl:

```
(CHARACTER C y
  (CHARWD R 0.5278)
  (CHARHT R 0.43055)
  (CHARDP R 0.194443)
  (CHARIC R 0.008)
```

lmr10.afm:

```
C 121 ; WX 527.77777 ; N y ; B 19 -205 508 431 ;
```

During the processing we check, compare and test for compatibility the metric data, especially the character widths, ligatures and kerning pairs crucial for typesetting taking in account T<sub>E</sub>X font limitations: A font contains at most 256 character codes, at most 15

different nonzero widths, 15 different nonzero depths, and 63 different nonzero italic corrections.

*Difference between  $\TeX$  and outline fonts.* The ‘native’  $\TeX$  fonts (pk/tfm) have no more than 256 characters, no character names, any comparison using  $\TeX$  is reduced on character sets defined by available encodings. On the other hand, the PostScript Adobe Type 1 fonts (pfb/afm) have glyph names, may contain many glyphs, but only 256 could be encoded, their encodings may be flexible, and all glyphs may be compared by names. Additionally, in OpenType (otf) many glyphs are encoded and available. Therefore, operations with an outline font could be independent off the  $\TeX$  limitations and all glyphs present in the font can be processed.

*Font tables and font specimens with  $\TeX$ .* To test a completeness of the font glyph set we start with font tables and font specimens. For this task testfont: CTAN:texmf/tex/plain/base/testfont.tex from the standard  $\TeX$  distribution, fonttabs: CTAN:texmf/tex/csplain/fonttabs.tex and OFS [4] developed by Petr Olšák could be recommended.

*METAFONT and bitmap fonts.* METAFONT generates the metric tfm files and a bitmap representation of glyphs for a selected device (mode). The shapes of bitmap fonts are represented as a bitmap (a matrix of pixels) in *any* resolution. Unfortunately, probably no reference resolution exists. We run METAFONT and dvips with modes available from texmf/metafont/misc/modes.mf

```
mf '\mode='$MOD';' input font.mf
dvips -mode $MOD -D $RES
```

and with the corresponding resolutions, for example

MOD	300	600	1200	2400	2602	5333
RES	cx	ljfour	ljfzzz	supre	proof	crs

to test fonts for all designed sizes and also for various (low, middle and high) resolutions. There is no direct correlation between correctness or incorrectness shapes in different design sizes or different resolutions (magnifications).

*Outline fonts.* Outline fonts, e.g. Adobe PostScript Type 1 and OpenType, are represented by their outline contour curves independent of resolution. Of course, rendering always depend on an output device for both outline and bitmat fonts. Good outline fonts should

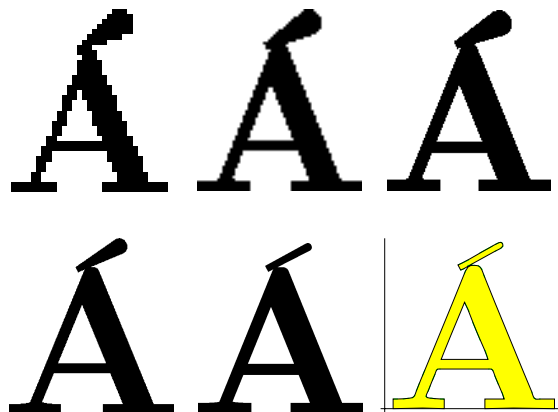


Figure 1: csbx10: Á at various resolutions.

be rendered properly elsewhere, for example, may be embedded in PDF document and are more flexible than bitmap fonts because “rerendering” of a bitmap font for any device may often cause the loss of quality. The aim of testing outlines is to verify *consistency* (font and glyph element are unified for all fonts of a font family, preserved for all design sizes) and *compatibility* of versions, changes may be only improvements or corrections of mistakes (not producing new mistakes).

#### PROOFPRINTING OF METAFONT FONTS

*csbx10: “Á”.* Our approach is not to prove correctness of METAFONT programs but to check their products — glyphs in the pk format. Because it is impossible to choose one resolution to verify we should test the glyphs for various resolutions. Fig. 1 shows tests of Á from the csbx10 font at following resolutions: 300, 600, 1200, 2602, and 5333 dpi. The last is the result of mftace (autotracing bitmaps). We can detect a bug in CSfonts [T $\text{\E}$ Xlive 2005] — serif upper case accents depend on resolution (mode). We have also to observe a consequence of such kind of bugs: the results of conversion to the outline font by autotracing high resolution bitmaps cannot be correct.

*cmbx9/cmbx10: “y”.* The next example illustrates behavior of “y” in bold CM fonts (see Fig. 2 with 600, 1200, 5333 dpi, and mftace) where a strange scrap is generated. This bug in CM is very small, its occurrence is more rare and more hidden (only in some sizes), it is not evident like the previous bug in CS. A similar effect can be observed for cmbx9, cmbx7, cmbx8 (see also [2]). Probably, weaker correlations suppress this



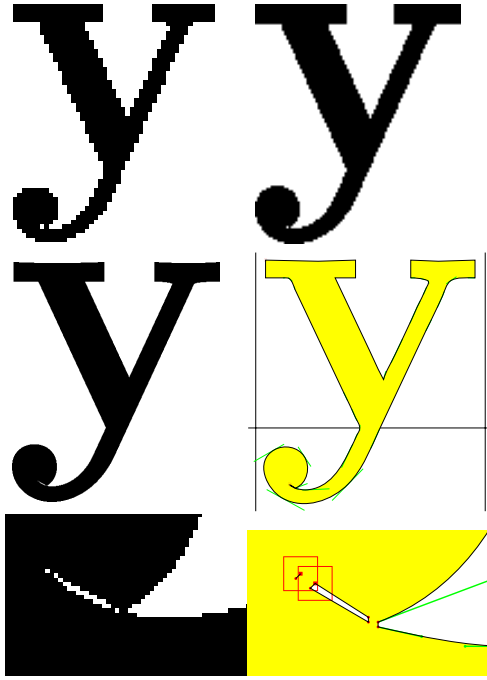


Figure 2: cmbx10: y.

defect for cmbx5, cmbx6, cmbx12, and cmbx17.

*The role of METAFONT today.* The METAFONT fonts may be simple or very complex, therefore debugging of a bitmapped glyph representation may be difficult. Although probably the users do not expect bitmap fonts in distributions today, METAFONT and METAPOST are still and will be a very important tool for font developers.

#### PROOFPRINTING OF OUTLINE FONTS

Figure 3 shows the proofing page produced from the LM sources directly by programs from the MetaType1 package [3]. The following pictures (in Fig. 4) demonstrate my own variants of proofsheets for screen and printer to recognize better the tiniest details (using zoom) invisible in the previous figure because of mutual overlapping such a small details, e.g. the control points and vectors. Figures 4 and 5 also open questions about an optimal approximation and a hinting strategy: To add the nodes at extrema or not, how to hint accents, etc.

FontForge [5], an open source font editor, developed by George Williams allows to read and generate various font formats, we can also use it for checking and analysis fonts, import of T<sub>E</sub>X font bitmaps, export glyph outline curves in eps for subsequent processing.

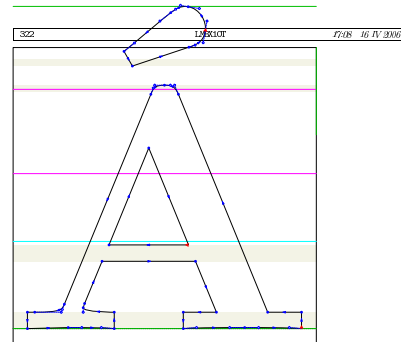


Figure 3: Proofprinting with MetaType1.

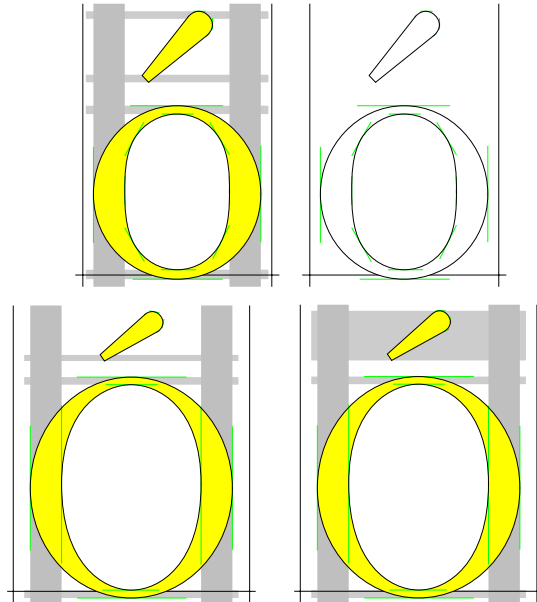


Figure 4: Proofprinting of outline fonts

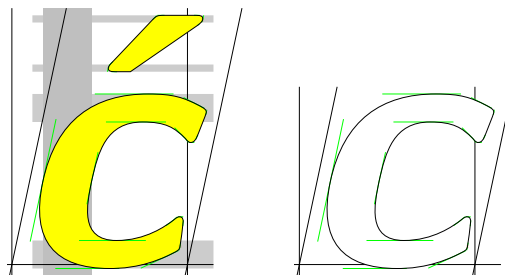


Figure 5: Extrema points and hints.

P.P.	AcAc	cmr10
P.P.	AcAc	ec-lmr10
P,P,	AdAd	cmr10
P,P,	AdAd	ec-lmr10
F.F.	AeAe	cmr10
F.F.	AeAe	ec-lmr10
F,F,	AoAo	cmr10
F,F,	AoAo	ec-lmr10

Figure 6: Kerning pairs in a visual form.

ajja]	ec-lmr10 <sub>0.99.3</sub>
ajja]	ec-lmr10 <sub>1.00</sub>
ajja]	ec-lmr10 <sub>0.99.3</sub>
ajja]	ec-lmr10 <sub>1.00</sub>
ff!f!	ec-lmr10 <sub>0.99.3</sub>
ff!f!	ec-lmr10 <sub>1.00</sub>
ff?f?	ec-lmr10 <sub>0.99.3</sub>
ff?f?	ec-lmr10 <sub>1.00</sub>

Figure 7: Kern changes in LM.

## COMPARISON OF METRIC DATA

Analyzing the metrics we detect (automatically) cases of agreement or disagreement of values of dimensions or kernings. In fig. 6 we present a small part of comparison between the CM and LM tfm files in T1(ec-lm) encoding. cmr10 and ec-lmr10 are compatible in presence of the kerns "P" : " , " . " and absence of the kerning shifts "F", "T", "V", "W", "Y" : " , " , " . ".

In ec-lmr10 new kerning pairs have been introduced: "A" : "c", "d", "e", "o".

Fig. 7 demonstrates the changes between two versions of LM.

```
\newlength{\bbox}\newlength{\cbox}%
\def\fbboxsep{0pt}\def\fbboxrule{0.1pt}
\def\pair#1#2{%
\settowidth{\bbox}{#1#2}%
\settowidth{\cbox}{\mbox{#1}\mbox{#2}}%
\addtolength{\bbox}{-\cbox}%
```

```
\fbbox{#1}\kern-0.2pt\kern\bbox\fbbox{#2}\fbbox{#1#2}
}
\pair{a}{j}
```

A short “ $\text{\TeX}$ ” macro pair calculates the shift between two “kerned” boxes.

## COMPARISON OF GLYPH IMAGES

We will demonstrate two techniques

- color mix with pdf $\text{\TeX}$  for visual scan
- outline comparison for outline fonts

*Color mix with pdf $\text{\TeX}$ .* Generating of comparative proofsheets using pdf $\text{\TeX}$  with mixing colors is possible for both bitmapped and outline fonts. Searching differences needs a subsequent human visual postprocessing. In our examples, combination of red and cyan produces pink in the intersection.

```
\usepackage{graphicx}
\def\Default{\pdfliteral{0 g 0 G}}
\pdfpageresources
{/ExtGState
<< /Luminosity
<< /Type /ExtGState /BM /Luminosity >>
>>}
```

```
\def\Acolor{1 0 0 rg}% red
\def\Bcolor{0 1 1 rg}% cyan
```

```
\renewcommand\C{\char#1}%
\Default
\fbbox{\makebox[0pt][l]%
{\pdfliteral{/Luminosity gs \Acolor}\Afont\C}%
{\pdfliteral{/Bcolor}\Bfont\C}}%
```

Fig. 8 demonstrates data from two samples:

*left* outline font lmbx10 (ver. 0.99.3) [cyan] v.s. bit-map font csbx10 [red]

*right* lmbx10 (ver. 0.99.3) [cyan] v.s. lmbx10 (ver. 1.00) [red] (both outline Type 1)

In a gray-scale printing the red color looks dark, the cyan is lighter, and the intersection is the lightest.

*Comparison of outline fonts.* Some elements of two outline fonts allow an automatical comparison. We can compare two fonts with a common glyph repertoire, e.g. to test two releases, alternatives, extensions or subsets of any font or to test two similar fonts for differences. We compare all the glyphs available in both fonts by their names automatically and detect the following differences:

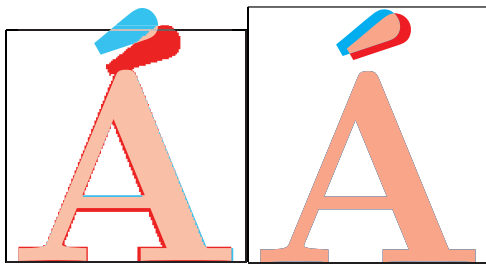


Figure 8: A color mix.

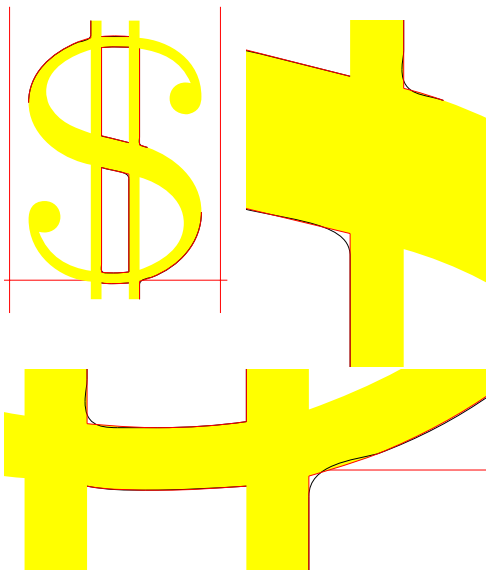


Figure 9: Improvement of outline approximation.

presence and absence of a glyph with a given name (it may signal a new glyph has appended, an old glyph has been removed; a glyph has been renamed, or sometimes, a glyph name may be invalid)

different glyph shapes (at least one segment is different),

different glyph widths (even after rounding to integer in the glyph coordinate space)

The contour curves of the first (older) font are black. The contour curves of the second (newer) font are red and the filled glyph area is yellow. The glyph box frames in the older font are blue.

The black and dark lines denote the older version in printing without colors.

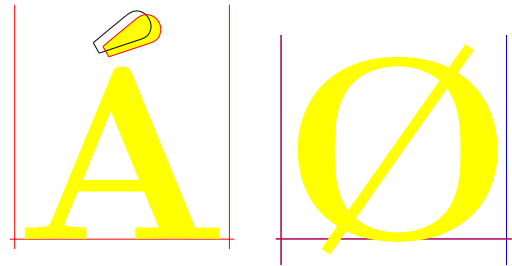


Figure 10: Modification and correction.

Figures 9 and 10 show differences between LM 0.99.3 and LM 1.00. In fig. 9 the dollaroldstyle from lmr10 has been improved (its conversion to outlines is better). Fig. 10 gives an information about modification of acute accent and correction of glyph width in the lmbx10 font.

## CONCLUSION

Techniques extending and complementary to existing testing tools in examples have been presented to help to font authors and maintainers in their time consuming and expensive work. The results of tests were or may be applied as warnings, bug reports or suggestions. Namely, they were or are used for verification of the Type 1 version of public Indic fonts [6] available from CTAN and for tests of the LM fonts.

## REFERENCES

- [1] Donald Knuth. METAFONT book.
- [2] Bogusław Jackowski, Janusz M. Nowacki, Piotr Strzelczyk. “Programming PostScript Type 1 Fonts Using MetaType1: Auditing, Enhancing, Creating. *Proceedings of the XIV EuroTeX 2003 conference*, pp. 151–157, Brest, France, 24–27 June 2003.
- [3] MetaType1 distribution: <ftp://bop.eps.gda.pl/pub/metatype1>.
- [4] Petr Olšák. Font management system OFS. <ftp://matf.feld.cvut.cz/pub/olsak/ofc>
- [5] George Williams. FontForge: an outline font editor, <http://fontforge.sourceforge.net>
- [6] CTAN:fonts/ps-type1/indic



# Typography – the art of the letter system

ISTVÁN RADÓ

## ABSTRACT

*The dialectic of the typographic print and the service of the ergonomic of reading; the damage caused by malformed prints resulting from their partial or complete negligence – in the context of European Roman letters.*

## EXTENDED ABSTRACT

Writing is conserving human thought in a systematic order. Printing is the generator of intellect, and occasionally a bearer of real knowledge, too. For over a thousand years letters have been the means of representation of the creative human intellect on papyrus, then paper, and for some 30 years now on the computer screen, too.

The term letter signifies the so-called Roman letters that emerged from ancient Greek and Roman culture. It was this letter type that had it as its function to be a porter of knowledge and culture in Europe, and thus it has served a worldwide technical civilisation emerging from European culture, with all its well-known blessings and curses.

Today 60% of humanity is illiterate, and half of the rest is functionally illiterate. If we keep halving, we find the percentage of those who do not read anything except the news and tabloids, therefore it is a mere 20% of the entire population of the earth that seeks knowledge from letters. And there is an ever-growing proportion within this 20% representing Chinese speakers! This latter fact provides a peculiar (and alarming) future to the Roman letter.

However, the documentation of the past, present and near future is still a task where European Roman letters and numbers dominate exclusively. The appearance and the point of this form of expression is the dialectic harmony of form and content, or an apparent and highly disturbing lack of it. It is easy to express the gist of a given idea in letters: unity of thought developing from contrapositions, opinion, reasoning, assertion, informative teaching, the obscuring of the point, plagiarism, incitement, forgery, lying etc.

A good, well-made and intelligent typography gives emphasis and rank to the thought it represents. A print type that is ostentatious, or one that neglects or consciously infringes a moderate proportionality is always a sign of the superficiality and false nature of the message itself. The blissful development of PCs regards as uniform and (perhaps unduly) degrades, too, all prints produced by repetitious use of the same font types and groups. The chaotic and disproportionate location of font size and spacing of lines are evidence of the damage made to the message, or its very worthlessness.



# Introduction into BibT<sub>E</sub>X style file Programming

---

■ BERND RAICHLE  
DANTE e.V.  
bernd.raichle@gmx.de

## ABSTRACT

*This tutorial shows how to generate bibT<sub>E</sub>X style file, and then how to modify it by hand when necessary. It also gives an introduction to the BibT<sub>E</sub>X style file language.*

## EXTENDED ABSTRACT

When you have to create a custom BibT<sub>E</sub>X style file, tools like “custom-bib/makebst” will simplify your life. Nonetheless it is often necessary to make small manual changes to the resulting or an existing .bst file to satisfy your needs.

This tutorial will enable BibT<sub>E</sub>X style file pro-

grammers and normal BibT<sub>E</sub>X users to get accustomed to the structure of a .bst file allowing them to make small changes to an existing BibT<sub>E</sub>X style file. To achieve this, we will start with a minimal .bst file showing the basic blocks. This minimal style file will be extended step by step introducing more and more .bst primitive functions.





# *pdfT<sub>E</sub>X – what was, is and will be*

---

MARTIN SCHRÖDER  
pdfT<sub>E</sub>X maintainer  
martin@oneiros.de

## ABSTRACT

*The talk is a review of key pdfT<sub>E</sub>X features and primitives presented on a timeline from the beginning, through present versions (1.30 and 1.40) till the features planned in the near future.*



# *L<sup>A</sup>T<sub>E</sub>X programming tutorial*

PÉTER SZABÓ

Budapest University of Technology and Economics,  
Department of Computer Science and Information Theory,  
H-1117 Hungary, Budapest, Magyar tudósok körútja 2.  
pts@fazekas.hu

## ABSTRACT

*This tutorial is a practical introduction to L<sup>A</sup>T<sub>E</sub>X programming: implementing new features (writing L<sup>A</sup>T<sub>E</sub>X packages), writing packages accepting options, changing existing features, finding out what commands to change, finding the file containing the definition of the command, overriding the definition, extending the definition, debugging, writing code independent of catcode changes, string processing and .aux file tricks.*

## FURTHER READING

*The Not So Short Introduction to L<sup>A</sup>T<sub>E</sub>X 2<sub>ε</sub>.* This is about using L<sup>A</sup>T<sub>E</sub>X for typesetting, not programming, but this is a good introduction to its syntax and main concepts. Translations available to several languages. <http://www.ctan.org/tex-archive/info/lshort/english/lshort.pdf>

*The T<sub>E</sub>Xbook.* Although it is about plain T<sub>E</sub>X, it explains some really advanced topics about T<sub>E</sub>X and its macro programming language, most of them being relevant to L<sup>A</sup>T<sub>E</sub>X, too. Paragraphs and exercises marked with single and double dangerous bends are especially recommended for thorough reading: these are the most authentic and in-depth explanations about how T<sub>E</sub>X works. Introductory exercise: try to download the T<sub>E</sub>Xbook from CTAN and compile it for yourself.

*The documentation of ε-T<sub>E</sub>X.* It documents some important new primitives. L<sup>A</sup>T<sub>E</sub>X now uses ε-T<sub>E</sub>X by

default, so these powerful primitives are available for the L<sup>A</sup>T<sub>E</sub>X programmer.

*The manual of pdfT<sub>E</sub>X.* It documents some important new primitives. This will help you understand how the pdf<sub>tex</sub> drivers of graphics.sty and hyperref.sty work. Compilation hint: download the manual folder with the file pdf<sub>tex</sub>-t.tex. Compile it with `texexec -pdf pdftex-t`. If the compilation falls to an infinite loop, abort it when pdfT<sub>E</sub>X finishes running.

*The developer documentation of the L<sup>A</sup>T<sub>E</sub>X 2<sub>ε</sub> kernel.* Download source from <ftp://ftp.ctan.org/pub/tex/macros/latex/base.zip>, compile .dtx files with L<sup>A</sup>T<sub>E</sub>X. It also contains the source code of the base document classes.

*The source and the developer documentation of graphics.sty, hyperref.sty and other advanced packages.* Fetch them from CTAN.



# Managing a math exercise database with L $\text{\TeX}$

PÉTER SZABÓ

Budapest University of Technology and Economics  
Dept. of Computer Science and Information Theory  
H-1117 Hungary, Budapest, Magyar tudósok körútja 2.  
pts@fazekas.hu

ANDRÁS HRASKÓ

Fazekas Mihály Fővárosi Gyakorló Ált. Isk. és Gimnázium  
Horváth Mihály tér 8.  
H-1082 Budapest, Hungary  
hraskoa@fazekas.hu

## ABSTRACT

*L $\text{\TeX}$  is a good tool for creating beautiful books, especially when the book contains a lot of math formulas. It is not rare that L $\text{\TeX}$  is used to typeset a view of a database, by generating L $\text{\TeX}$  source from the database text, possibly using XML as an intermediate format. Some L $\text{\TeX}$  packages and formats support reading XML data directly.*

*In the matbook project we have created a database of math exercises for special class secondary school students, and as well solutions and instructions for teachers. The data is organized in a tree structure of custom L $\text{\TeX}$  environments in .tex source files. L $\text{\TeX}$  reads these data files several times for generating the books. CVS is used for data replication and concurrent co-authoring. We are planning to convert to use a L $\text{\TeX}$ -to-HTML translator to publish the database on the web.*

*This paper presents the simple software architecture of the matbook project and the design decisions we made concerning software and workflow, and it also compares matbook with other approaches such as big content management systems and L $\text{\TeX}$ -enabled wikis.*

## NON-STANDARD USE OF L $\text{\TeX}$

The original purpose of L $\text{\TeX}$  (and L $\text{\TeX}$ ) is typesetting beautiful books, journals and other printed material.

Novel uses include preparing slides for talks, developing software and its documentation together (e.g. web and ltxdoc), typesetting math formulas (e.g. Texvc [11]), typesetting printed and on-line HTML documentation together, rearranging PDF pages (pdfL $\text{\TeX}$  with pdfpages.sty) and typesetting text generated from databases or other markup formats.

In the matbook project we use L $\text{\TeX}$  to read a database of math exercises (in several passes), and typeset the material to books for students and teachers. This paper presents the software architecture and some implementation details of the matbook project, and it is also a case study of integrating excellent free software tools for low-budget publishing.

## PROJECT GOALS AND PRODUCTS

The *Fazekas Mihály Secondary Grammar School* of Budapest [1] has been launching special mathematics classes for several decades, and is proud of its students winning national and international student competitions, and later becoming appreciated mathematicians. E.g. László Lovász, the well-known Hungarian mathematician graduated in Fazekas in 1966.

Good mathematicians have good problem solving skills, and this skill can be best developed by solving problems and exercises. It is the responsibility of the teacher to choose the exercises for the students which best fit their learning curve. Talented students in a special math class need special attention. A lot of exercises and didactic experience have accumulated in Fazekas over the last few decades, and we have decided to publish this in printed form in Hungary, and we are also planning to provide a web interface where all material is available. Thus matbook was born.

We are compiling a comprehensive exercise database (which also includes solutions, didactic advices, exercise lists for lessons and metadata for more accurate searching). Students and teachers in Fazekas are both working on extending this database, and we are developing software that would present this database to its audience. We are planning to publish exercise books (for students) and teachers' guides. If students buy the exercise books, teachers can give homework assignments from those books. (Of course, teachers will assign exercises whose solutions cannot be found in the exercise book.)

We are also planning to provide a web interface on which visitors can browse and view exercises, solutions etc., they can do a full text search, and they

can also search for exercises in a given topic (specified using a set of predefined keywords). We already have a web interface for a comprehensive database of Hungarian secondary school math contest problems (which stores text in  $\text{\LaTeX}$  format, and converts it to HTML using TTH [2]), and we'd like unify this with the matbook database.

## DATABASE STRUCTURE

The database consists of

- *exercises* for the students;
- *hints* and *solutions* corresponding to the exercises, for the students;
- solutions for the teachers only;
- *remarks* and *didactic advice* corresponding to the exercises, for the teachers;
- a hierarchic taxonomy of *keywords* covering the fields of mathematics (e.g. prime numbers, trigonometry);
- association between exercises and keywords;
- organization of exercises to *chapters* and *volumes*;
- chapter and volume introduction text;
- ordered *exercise lists* prepared for obligatory and facultative lessons.
- *figures* referred to in the text, in EPS format

## SOFTWARE COMPONENTS

- *volume typesetter*: a set of  $\text{\LaTeX}$  macros to read the database in multiple passes, and typeset the book volumes;
- *indexer*: generates the keyword index at the end of the volumes (similar to *makeindex*);
- *web user interface*: with browse, view and search functionality;
- *consistence validator*: checks whether database files conform to the specifications.

Existing free software used: standard tools in a  $\text{\TeX}$  distribution, the lmodern font family [3], GNU Ghostscript, sam2p [4], ImageMagick, CVS, Perl, the new magyar.ldf (part of [5]), husort.pl (Hungarian index processor, part of [5]), stuki.sty (structogram figure generator [6]).

We work in a Linux-Windows mixed environment, so it was our aim that all components except for the server part of the web user interface should run on both UNIX and Win32.  $\text{\TeX}$  tools we need are available on both systems. We decided to implement the indexer and the consistence validator as command-line Perl applications so it would be easy to port them across systems.

## DATABASE LAYOUT

We chose storing our data in structured text files rather than using a relational database, because it is easier to change the schema later, and we don't have to develop a custom user interface for data editing. XML is a good and widely supported structured text data model and syntax, but we prefer a format which is quick to type and easy to review for humans. YAML [7] is such a format. We finally chose the XML data model (for interoperability with other software), but a  $\text{\LaTeX}$ -compatible syntax (for easy typing), which can be converted to XML without loss when needed.

As a master text markup format, we quickly rejected XHTML + CSS + MathML, mostly because it is tiresome to type a document in this format; it is not possible to archive a rendered version of an XHTML text in a scalable way; it is not possible to specify typesetting hints (such as penalties); and MathML is not powerful enough: it is not possible to type the right, textual side of  $\backslash\text{cases}$  in MathML; MathML still lacks some symbols. Moreover, with current browsers it is not possible to ensure acceptable visual quality: browsers render the same document differently, MathML support usually doesn't come out of the box, browser MathML fonts lack important symbols, browsers cannot hyphenate long words automatically, the visual output depends on the installed fonts and the browser window size (which the author of the text cannot control), browser cannot break the line in the middle of a MathML formula etc.

We could have adopted a safe and easy to type markup format, similar to MediaWiki's WikiText format [8] or  $\text{\SäferTeX}$  [9]. The MediaWiki software implements the text rendering engine of Wikipedia [10], and it lets authors insert math formulas in a subset of (AMS) $\text{\LaTeX}$  syntax. When the page is rendered, these formulas are interpreted and converted to images or MathML formulas by Texvc [11]. We have rejected MediaWiki because – similarly to XHTML

– it doesn't give the author enough power to ensure perfect visual output quality. We did not use ŠäferT<sub>E</sub>X because its source code was not available, and it was not mature enough.

We could have invented our own markup format. Doing this would have required us not only to invent an excellent format, but to write a renderer (to both PDF and HTML), and document the format thoroughly, including tutorials and examples. This option was not feasible in our project.

Thus we have chosen a restricted subset of L<sup>A</sup>T<sub>E</sub>X as a markup format. Its advantages are: it has been available for a long time, the basic command set is well-documented, it gives the text author sufficient control over the visual quality of the output, and there are lot of fonts and packages we can use. We had to impose restrictions in order to keep our format convertible (primarily to XHTML+CSS+MathML). The most important restrictions on the document text are: it is forbidden to load packages or other files, define or change macros, use conditionals or other programming features, change catcodes, use the character " in the input (but the "proper quotes" must be used), use conditionals, insert figures with `\includegraphics` (we provide a more restricted command instead).

Once we settled on L<sup>A</sup>T<sub>E</sub>X as a text markup format, it was straightforward to use the same syntax for structuring the data, so that our database text files won't contain two alternating formats, and they can be syntax-highlighted or otherwise processed in text editors easily. However, plain L<sup>A</sup>T<sub>E</sub>X is not suitable for structuring. For example, it is not obvious to deduce where chapter "First" ends in this L<sup>A</sup>T<sub>E</sub>X source, without knowing the meaning and depth of `\section`:

```
\chapter{First}    \emph{First} content.
\section{Inside}   \emph{Inside} content.
\chapter{Second}\label{2nd} % dummy
                  \emph{Second} content.
```

The XHTML representation (using `<H1>` for chapter titles and `<H2>` for section titles) suffers from the same limitation.

Our data format solves the problem by specifying structure using custom L<sup>A</sup>T<sub>E</sub>X environments. The example above looks like this:

```
\begin{mchapter}{title={First}}
  \emph{First} content.
\begin{msection}{title={Inside}}
  \emph{Inside} content.
```

```
\end{msection}
\end{mchapter}
\begin{mchapter}{title={Second},id={2nd}}
  % dummy
  \emph{Second} content.
\end{mchapter}
```

When converted to XML, it becomes:

```
<mchapter title="First">
  \emph{First} content.
  <msection title="Inside">
    \emph{Inside} content.
  </msection>
</mchapter>
<mchapter title="Second" id="2nd">
  <!-- dummy -->
  \emph{Second} content.
</mchapter>
```

Please note that `\emph` is not converted, because it is part of the text markup, and not part of the structure.

Thus there is a simple mapping between XML and our data format:

- L<sup>A</sup>T<sub>E</sub>X environment with attributes (i.e. "key = value" pairs) ↔ XML tag with attributes, properly escaped
- T<sub>E</sub>X comment ↔ XML comment
- other L<sup>A</sup>T<sub>E</sub>X text ↔ XML text (PCDATA)

This direct mapping makes it possible to use XML tools on our database. For example, we can use XSLT to do structural transformations on the XML, and we can use DTD or XML Schema validators to validate our database.

## DATABASE FOLDERS AND FILES

The database is spread into several small text files in several folders. The files read each other (using `\input`). The points where the data must be split and the naming conventions for the files and folders are strictly regulated.

The database is replicated on each co-worker's machine, using the CVS [12] revision control system. People can work offline, and commit their changes back to the repository on the server a few times a day. Server failures and network connection slowdowns don't affect working hours seriously. CVS is smart enough to merge concurrent but independent changes of text files, and it enforces human interaction when a conflict occurs, so there is no danger of accidentally overwriting somebody else's changes. CVS also keeps

old versions too, so accidentally deleted text can be recovered any time later. (CVS merges files line-by-line, which complicates concurrent editing of binary files hard – but this limitation doesn’t affect our project since we mostly use text files.) Subversion (= SVN [13]) is a newer and more advanced revision control system, and it won’t hard to migrate from CVS when those advanced features are needed. Both CVS and Subversion have clients on multiple platforms, including UNIX and Win32. We use the standard cvs client on Linux, and TortoiseCVS on Windows.

The reason why the database is split to multiple files is that it is easier to transfer changes of smaller files in CVS, and it is also easier for humans to edit a few small files concurrently than to edit one large file. Usually multiple people are modifying the database at the same time, but most of the time they work in their own files, so no conflict occurs. Using multiple folders makes it easier to select the correct file for opening.

The file and folder layout also follows the L<sup>A</sup>T<sub>E</sub>X compilation process. Compilation always starts in the root folder (of the CVS tree). For example, to compile the first volume of Algebra, one runs “`latex volume_a_i`”, which starts processing the file `volume_a_i.tex`. All other files belonging to this volume are reside in the folder `chs_a_i` and its subfolders. Files `\input` are thus specified relative to the root folder, thus adding `../` is not necessary when referring to local `.sty` files. It is also convenient that all temporary and output files go to the root folder, thus subfolders are not changed during the compilation process.

### L<sup>A</sup>T<sub>E</sub>X TRICKS

In this section we present some problems we faced when typesetting with L<sup>A</sup>T<sub>E</sub>X; solutions included.

**UNIFIED LABELS** This is a feature that makes it possible to refer to a `\label` defined in another volume. It is accomplished by reading `\newlabel` commands from the other `.aux` files, and adding them with both the label text and page numbers prefixed with the other volume name.

**VOLUME SPLIT** Since teachers’ guides can be several hundred pages long, it might be necessary to split them to multiple volumes. If it is so, the editor promotes some chapter boundaries are to volume boundaries by adding the appropriate command to the source of the main `.tex` file. We chose the most

portable ways to typeset these subvolumes: all subvolumes are separate L<sup>A</sup>T<sub>E</sub>X documents, which `\input` the main `.tex` file in a mode in which the `\shipout` of the unnecessary pages is cancelled.

The advantage of this method is that it doesn’t require external tools (such as `pselect`), it works for both PostScript and PDF, and it can be run from a regular L<sup>A</sup>T<sub>E</sub>X IDE. It’s disadvantage is its slowness. An alternative approach for PostScript would be generating and running a `pselect` command line for each volume. And for PDF, an alternative approach is selecting the appropriate pages from the main volume using `pdfpages.sty`.

Splitting the volume into real subvolumes (so that compiling a subvolume doesn’t read the other subvolumes’ L<sup>A</sup>T<sub>E</sub>X source) would not work, because it would render page numbers, the bibliography and inter-subvolume references wrong, which would need additional programming to get right.

**FUZZY KEYWORD NAMES** When specifying the list keywords associated to an exercise, it is big burden to specify a long keyword precisely (with spaces, punctuation etc.). In order to solve this, a Perl script generates keyword aliases, e.g. the first word of a keyword will become an alias to the keyword if it is not ambiguous.

**STRING PROCESSING** Another idea in fuzzy keyword name matching to match keywords after stripping spaces, punctuation, upper case and accents. This stripping had to be implemented in L<sup>A</sup>T<sub>E</sub>X, too. Since T<sub>E</sub>X doesn’t have string processing primitives, we have to implement them using macros. Here is a mix of a macro definitions that shows the most important string processing tricks:

```
\def\stripit#1>{\def\empty{} \def\space{ }
\def\rmonestar#1{\ifx#1\hfuzz\empty\else
\if*\string#1\else#1\fi
\expandafter\rmonestar\fi}
\begingroup\lccode‘!‘ \lowercase{\endgroup
\def\oonespace#1 {\ifx\hfuzz#1\empty\else
#1!\expandafter\oonespace\fi}}
\def\rmstars{%
\afterassignment\rmstarsb\def\M}
\def\rmstarsb{%
\edef\M{\expandafter\stripit\meaning\M
\space\hfuzz\space}
\edef\M{\expandafter\oonespace\M}
\edef\M{\expandafter\rmonestar\M\hfuzz}}
```



The macro `\rmstars` above removes all stars (\*) from a string. The string is specified as an argument in braces, and the result – without the stars and all tokens having catcode 12 – it is put into the macro `\M`. Example invocation: “`\rmstars{a * B**cd} \show\M`”.

A rough outline of its operation: `\meaning` converts tokens to catcode 12, except for spaces, which are converted to catcode 10. Then `\oonospace` iterates over all spaces and converts them to catcode 12, too. Finally `\rmonestar` iterates over the tokens and removes all stars. Almost beautiful.

Read more about the primitives involved here in The T<sub>E</sub>Xbook [14].

ONE OR MORE SOLUTIONS? If there is one solution for an exercise, it should be prefixed with “Solution” (and not “Solution 1”). If there are more solution, each of them should have a number, “Solution 1”, “Solution 2” etc. By the time of emitting the 1st solution, we don’t have the information whether there are more. How do we typeset it properly?

To solve problems like that, it is a common trick to use the `\label`–`\ref` mechanism. We emit `\label{exercise42-sol2}` at “Solution 2”, and the next the the document is recompiled, at “Solution 1” we check for the presence of this label, e.g. with

```
\@ifundefined{r@exercise42-sol2}{...}{...}
```

BIBLIOGRAPHY THREE TIMES When a `\cite` command with a new a target is added to the document, it is necessary to run L<sup>A</sup>T<sub>E</sub>X three times: `latex doc; bibtex doc; latex doc; latex doc`. The 1st run of L<sup>A</sup>T<sub>E</sub>X records the `\citation` command to the `.aux` file. The BibT<sub>E</sub>X run generates the `.bbl` file. The 2nd run of L<sup>A</sup>T<sub>E</sub>X inserts the new `.bbl` file to the document, and it also records the `\bibcite` command to the `.aux` file indicating the new number to be displayed by the `\cite` command. The last, 3rd run of L<sup>A</sup>T<sub>E</sub>X indicates makes `\cite` emit that number.

We speed this up by parsing the `.bbl` file at the beginning (before the first `\cite`), so the 3rd run of L<sup>A</sup>T<sub>E</sub>X is not necessary.

## CONCLUSION AND FUTURE WORK

The matbook project demonstrates not only the power, openness and flexibility of L<sup>A</sup>T<sub>E</sub>X, but it is also an example of low-budget publishing using free software and a little scripting. matbook is also free software.

Our most important future goals are completing the exercise database and implementing the missing software components: a thorough consistence generator and the web user interface.

## REFERENCES

- [1] Fazekas Mihály Secondary Grammar School of Budapest. <http://www.fazekas.hu/>
- [2] TTH: the T<sub>E</sub>X to HTML translator. <http://hutchinson.belmont.ma.us/tth/>
- [3] Bogusław Jackowski and Janusz M. Nowacki. *Latin Modern fonts: how less means more*. <http://www.dante.de/dante/events/eurotex/papers/TUT09.pdf>, 2005.
- [4] Péter Szabó. *Inserting figures into T<sub>E</sub>X documents*. In proceedings to EuroBachT<sub>E</sub>X 2003.
- [5] Péter Szabó. *Implementation tricks in the Hungarian Babel module*. In proc. to TUG 2004.
- [6] Károly Lőrentey. *stuki.sty: Structograms in L<sup>A</sup>T<sub>E</sub>X*. <http://lorentey.hu/project/stuki.html.en>
- [7] YAML: machine parsable data serialization format. <http://www.yaml.net/>
- [8] WikiText: wiki markup language. <http://en.wikipedia.org/wiki/Wikitext>
- [9] Frank Schäfer. *ŠäferT<sub>E</sub>X: Source Code Esthetics for Automated Typesetting*. In proc. to TUG 2004.
- [10] Wikipedia: the free encyclopedia that anyone can edit. <http://en.wikipedia.org/>
- [11] Texvc: T<sub>E</sub>X validator and converter. <http://en.wikipedia.org/wiki/Texvc>
- [12] Karl Fogel and Moshe Bar. *Open Source Development with CVS*. 3rd Edition. O’Reilly, 2003.
- [13] Ben Collins-Sussman et al. *Version Control with Subversion*. O’Reilly, 2004.
- [14] Donald E. Knuth. *The T<sub>E</sub>Xbook*. Addison–Wesley, 1984.



# Would Aldus Manutius have used T<sub>E</sub>X?

ANDRÁS VIRÁGVÖLGYI

historian, designer

www.sequens.hu



Aldus' emblem, the dolphin curling around an anchor

## ABSTRACT

*Traditional Typography · Golden Mean · Linear Reading · The Siege of Constantinople · Aldus Manutius · His Ambitious Publishing Program · The Most Beautiful Book of the World · Festina Lente · Would Aldus have used T<sub>E</sub>X?*

It might be of interest to the participants attending the 16<sup>th</sup> EuroT<sub>E</sub>X conference to hear about the golden age of traditional typography in the 16<sup>th</sup> century and get to know one of the most famous book publisher of all times, Aldus Manutius of Venice.

The word typography is a compound of Greek elements. *Typos* means engraved illustration and *grapho* means to write.<sup>1</sup> *Scrivere sine penne* (to write without a pen) as they called it in the 16<sup>th</sup> century. Typography uses predefined shapes, letters, texts and sometimes illustrations as well. As final output it creates the printed page which represents another level of visual quality. After the new types of books of the 16<sup>th</sup> century left behind the poor readability and ponderousness of old codices, typographers who at the time were the printers themselves, followed the rules of the retrospectively labeled 'traditional typography.' It was a combined outcome of renaissance modern style and the lasting effect of medieval lettering. 'Traditional typography' determined the encounters of many generations with letters and enjoyed a consciously or unconsciously accepted status in the eyes of readers throughout the centuries.<sup>2</sup>

The design of books according to the traditional style starts with choosing the right proportions for page margins. Setting the inside, outside, top and bottom margins will give us the type area. The height, width and carefully chosen position of the type area influences overall proportions and balance of the book. Not everyone and not every workshop took the time and the trouble to carry out this kind of precision work and fine tuning. As time went by well tested

recipes started to form, incorporating certain elements of antique architecture and epigraphy. The *golden mean* for example, very popular among renaissance artists also found its way to book design. It stated that the smaller portion should compare to the bigger as the bigger compares to the whole, like in the case of the series 3 : 5 : 8. The relevance of these proportions could be observed in nature or in antique architecture. Typographers like sculptors, painters or architects before them, also wanted to take advantage of this noble rule when designing letters, margins or title pages.

After they made their decision about margin sizes, they could begin to compose the text. Information was organized according to a tight hierarchy, starting with the biggest type size used on title pages through the opening pages of chapters and the type of the body text until the smaller fonts of marginal notes and indices. The rhythm was provided by the repeating order of chapter headings, subheads and paragraphs all indicated by traditional typographic methods.

Further decisive features of the traditional page spread were symmetry and static arrangement. This structure encouraged linear reading so readers of the 16<sup>th</sup> century did not feel at all that they missed something. They read a book thoroughly from cover to cover. The list of recommended books was short and they had the time to periodically re-read some of them. Not to mention the Bible itself, which was studied and re-studied by many people every year.



Printers  
(16<sup>th</sup> century woodcut)

<sup>1</sup>Péter Virágvolgyi: *The Art of Typography*. [Osiris Handbooks] Budapest, 2002, Osiris Kiadó.

<sup>2</sup>Suzanne West: *Working with Style. Traditional and Modern Approaches of Page Design*. Budapest, 1998, UR Könyvkiadó. p. 54.

Renaissance typography gradually filled the space created by the newly invented art of printing. New and substantial methods of arranging information profoundly affected reception and the way of thinking of 16<sup>th</sup> century people as the digital world influences our approaches nowadays.<sup>3</sup>

### ALDUS MANUTIUS

In 1453 the Turks occupied Constantinople. Around the same time Gutenberg was preparing his 42-line Bible for printing. The taking of Constantinople not only coincided with the beginning of European printing but had an indirect but important effect on it as well. During Byzantine times many Greek schools prospered on the coast of the Bosphorus strait. When the Turkish siege started these Greek scholars left their schools and fled to Italy. They had an extensive knowledge of classic authors of ancient times. So their relocation turned Venice into a true center of classical sciences and research. There was something up in the air in that city, the situation just needed an entrepreneur to take advantage of it.

Aldus Manutius was a modest language instructor at the time, on the side of famous humanists like Pico della Mirandola. He taught them Greek and Latin and surely felt the need to create better editions to be used in education or for other humanists working with the same texts. Understanding the significance of all those Greek scholars arriving to Venice he moved to the city to try his luck. At first he was employed by a Venetian book merchant and printer called Asola. He was later named the leader of Asola's workshop and possibly in relation with this the owner's daughter be-



Aldus Manutius

came Aldus' wife. This was not an unusual step though, as it was Asola's best way to ensure that his talented employee would carry on with his thriving workshop. Aldus on the other hand had navigated himself to a position from which point he could set about executing his monumental plans. To secure a notable intellectual background he contacted Latin speaking humanists as well as the aforementioned Greek emigrants. He could frequently consult Greek scholars of Constantinople and hired scribes of the island of Crete as Greek typesetters and proofreaders.

The realization of his publishing program started in 1495. The result was the creation of probably the most beautiful and effectual books in the history of printing. Greek authors were published at first. The six volume Aristotle should be mentioned above all, which used a Greek cursive type modeling those scribes' original handwriting, and appeared in print between 1495 and 1498, causing a big sensation among European humanists. This initial success was soon followed by editions of Sophocles, Plato and Thucydides. After them came other classical authors writing in Latin: Virgil, Horace and Ovid. The size of each run were around a thousand copies. A surprisingly modern publishing policy governed these early critical editions. Aldus being a major representative of humanist approach of classical literature, preferred to get rid of all the medieval commentaries and foster the reading of original texts in their original languages. So it was up to the reader to make up his or her own variant, and to interact freely with ancient authors. To help learn the classical languages – since originally he was a teacher – he also published quality textbooks and dictionaries.

The most beautiful book of the world, as book historians like to call it, was published in 1499, four years after the initiation of Aldus' ambitious program. Francesco Colonna's allegoric and mythic poem, the *Hypnerotomachia Poliphili*, by 170 excellent renaissance wood-cut illustrations, fully harmonizing type design and margin proportions, received a form which still evokes the admiration of today's typographers and bibliophile book collectors all over the world. Surprisingly enough, though the specialty of *Hypnerotomachia Poliphili* was certainly noticed by contemporaries (a pirate edition came out soon after in Lyon), it was never reprinted in Venice within the next hundred years.

<sup>3</sup> Leah S. Marcus: *The Silence of the Archive and the Noise of Cyberspace*. In: *The Renaissance Computer. Knowledge and technology in the first age of print*. Eds. Neil Rhodes és Jonathan Sawday. London, 2000, Routledge. p. 22.

POLIPHILLO INCOMINCIA IL SECONDO LIBRO DI LA SVA HYPNER OTOMACHIA. NEL QUALE POLIA ET LVI DISER. TABONDI, IN QUALE MODO ET VARIO CASO NARRANO INTERCALARIAMENTI IL SVO INAMORAMENTO.

NARRA QVIVI LA DIVA POLIA LA NOBILE ET ANTIQVA ORIGINE SVA. ET COMO PER LI PREDECESSORISVITRI VISIO FVE EDIFICATO. ET DI QVEL LA GENTE LELIA ORIVNDA. ET PER QVALE MODO DISA VEDVTA ET INSCIA DISCONCIAMENTE SE INAMOROE DILE IL SVO DILECTO POLIPHILLO.



**M**E MIE DEBILE VOCE TALEOGRATIOE & diue Nymphie abfone peruenirano & incocine alla uoftra benigna audietia, quale laterifica rauceitate del urinate Efacho al fuo ue canto dela piageuole Philomela. Nondi meno uolendo io cum tuti gli mei exili conuen del intellecto, & cum la mia paucula fufficitia di fuffitare alle uofre piaceuole petitione, non rifaro al potere. L'eguale femota qualique hefitatione epie piu che fi congruerrebbe altronde, dignamente meritanio piu uertimo fluuto di eloquentia, cum troppo piu rotunda elegantia & cum piu exornata politura di pronuntiato, che in me per alcuno pacto non fi troua, di cofeguire il fuo gratiofo affecto. Ma a uui Celibe Nymphie & adme aliquato, quantiche & confufa & incomptante fringuliete haro in qualche pertiuncula gratificato affai. Quando uoluntatofa & diuota a gli defti uoftri & pofulato me prefato piu prefco cum lamino os mediocre prompto humile parendo, che cum enucleata terfa, & uenufa eloquentia placido. La pifica dunque & ueterima geneologia, & profapia, & il fatale mio amore garrulando ordire. Onde gia efendo nel uofro uenerando conuentuale confpecto, & uederme fertile & ieuua di eloquio & ad tanto prefite & di uo ceto di uui O Nymphie fedule famularie dal accefo cupidine. Et itanto benigno & delectuole & facro fito, di linere aure & florigeri fpiramini affiato. Io acconciamente compulfa di affumere uno uenerabile aulio, & tranquillo timore di dire. Dunque auante il tuto uenia date, o belliffime & beatiffime Nymphie a questo mio blacterare & agli femelli & terregeni, & puiffilli Conati, fit aduene che in alcuna parte io incautamente

A



POLIPHILLO QVIVI NARRA, CHE GLI PARVE ANCORA DI DORMIRE ET ALTRONDE IN SOMNO RITROVARSE IN VNA CONVALLE, LA QVALE NEL FINEERA SERATA DE VNAMIRABILE CLAVSURA CVM VNA PORTENTOSA PYRAMIDE, DE ADMIRATIONE DIGNA, ET VNO EXCELSO OBELISCO DE SOPRA. LA QVALE CVM DILIGENTIA ET PIACERE SVBTILMENTE LA CONSIDEROE.



**A** SPAVENTEVOLE SILVA ET CONSTITATO Nemoreufo & gli primi altri lochi per el dolce fomno che fe hauea per le felfe & profermate mebre difufio rediti, me ritrouai di nouo in uno piu delectabile fito affai piu che el precedente. El quale non era demon ti horridi, & crepidinole rupe intorniato, ne falcato di ftrumofi iugi. Ma compofitamente de grate montagniole di non troppo alicia. Siluofe di giovani queuoli, di rebuni, fraxini & Carpinii, & di frondofi Eculi, & llice, & di teneri Coryli, & di Alni & di Tili, & di Opio, & de infrutuofo Oleaftri, difpofiti fecondo la fpecto de gli arboriferi Colli. Et giu al piano erano grate filuole di altri filuatici

Pages from the  
*Hypnerotomachia*  
*Poliphili*, one  
of the most beautiful  
book of the world

Aldus' achievements deserve acknowledgement already, but his most remarkable innovation is here to follow. Production of books in the 16<sup>th</sup> century ever becoming cheaper and faster, workshops could start counting on bigger audiences for them. He was first to realize that instead the large format books printed before, readers who preferred solitary reading needed portable or pocket-size editions – as we call them today. Fifty years passed after the appearance of Gutenberg's enormous, two column 42-line Bible, when the Aldus Officina in Venice, leaving the old codex format behind, started to produce the incredibly popular one column pocket size editions of the classics.

Soon the *Aldi Neacademia*, a distinguished group of scholars (men of letters) was formed. By this time the workshop's Greek and Latin consultants had daily meetings to decide about the titles to be published. While the medieval scholar accumulated, the renaissance humanist judged the old manuscripts: they reconsidered the classic authors known in medieval times and tried to acquire previously unknown works by research or purchase. This was the way *editio princeps* books (first printed versions of the classics) got published, considered rare gems by later collectors. However they did not refuse to publish eminent medieval authors either, we can find Dante and Petrarch in Aldine editions. The farseeing editorial policy had its

fruitful result, the carefully selected titles were appreciated all over Europe.<sup>4</sup>

Characteristic to the success of the pocket-size editions is the large number of imitators especially in France. Simone de Colines launched a similar series in Paris, while Aldus had to defend himself against the pirate editions of the Lyon workshops by issuing a public protest letter. In it he enumerates the errors made by the printers of Lyon, so that the original could be easily discerned from the fake. The antecedent of the pirate editions of course was the appearance and success of Aldine publications in the French market. Classical Roman culture crossed the Alps in the form of these beautiful books, and as Beatus Rhenanus, the biographer of Erasmus put it: "northern barbarians" could now learn Latin and Greek while saving a long journey to Italy.

As an Italian Aldus was a faithful supporter of antiqua letters. In the beginning he used fonts strikingly similar to those used by the French Nicholas Jenson, also working in Venice. In 1501 he commissioned the talented letter designer, Francesco Griffo (1450–1518) to create a new Latin typeface. This was the first antiqua typeface which used capitals somewhat smaller than the ascenders of lower case letters. In the next year Griffo designed the famous cursive or *italic* type as it was later called, honoring the country of origin. Based on the so-called *cancelleresca*, it became the printed version of humanistic handwriting. Today we use italic to

<sup>4</sup> Martin Lowry: *The World of Aldus Manutius*. Oxford, 1979, Oxford University Press.

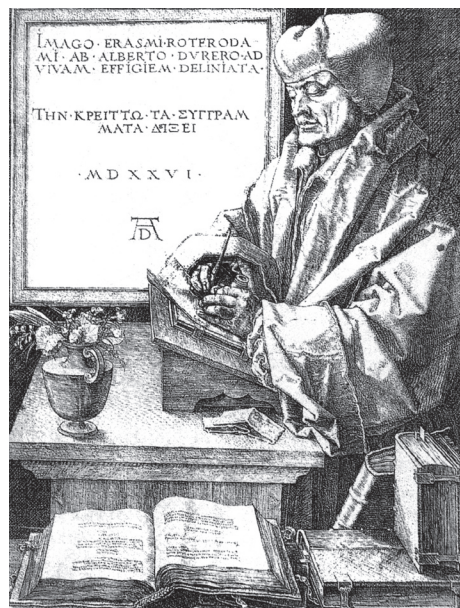


put emphasis on words or passages of text, but Aldus actually used it for typesetting whole books. Italics started their career in his workshop, exercising a great deal of influence on 16<sup>th</sup> century typography and fostering the victory of antiqua type all across Europe.

Page numbering also spread with Aldus' books, who was among the first to recognize its practical importance. Page numbers, apart from assisting the work of bookbinders, made much easier for readers to refer to a given section within a book.

The Aldus workshop had its heyday at the turn of the 15<sup>th</sup> and 16<sup>th</sup> centuries. In the Juvenile edition of 1501, Aldus set forth his philosophy and goals, many of whom were already achieved by releasing the series of elegantly designed and carefully edited classic publications. His logo, the dolphin with the anchor, appeared in his books from 1502. In the eyes of booklovers this mark represents utmost excellence considering both content and form. The dolphin curls around the anchor and emphasizes the classic saying which goes with it: *Festina lente* [gr. *Speude bardeos*], ie. to hurry slowly. It was Erasmus<sup>5</sup> who made this saying popular. Originally it comes from a play called *The Knights* by Aristophanes: "Speude takheos" – to hurry up quickly. In its reversed form the saying has several meanings. According to Erasmus this stoic statement should prevent princes from acting in the heat of the moment, to avoid swift and arbitrary decisions. Fabius Maximus is mentioned as the best example for slow diligence, who continuously weakened the invading army led by Hannibal employing his distressing technique. That is why he was labeled the postponer (*cunctator*) by Roman politicians. Supposedly emperor Augustus and Vespasian also liked the saying *Festina lente*. The dolphin curling around the anchor appears on one of the coins issued by Vespasian.

Posterity esteems Aldus' publishing activity highly. Jacob Burckhard (1818–1897), one of the earliest and most famous researchers of the renaissance highlighted his importance. Several contemporary memorial records the exceptional popularity of Aldine books. The following quotation comes from a letter by the German humanist Heinrich Glareanus to Ulrich Zwingli, dated from 19<sup>th</sup> of October, 1516: "I cannot miss to mention, that Wolfgang Lachner,



Portrait of Erasmus of Rotterdam by Albrecht Dürer (1526)

our Frobenus' father-in-law ordered a wagonful of classics from Venice, the best of Aldus' publications. If you would like to have some of them, let me know quickly and send cash. Because as soon as a similar shipment arrives, there are already thirty people surrounding it and keep asking "how much is it?" then start to fight over it. Passion rapidly becomes true fury and often seizes men who cannot even understand them."<sup>6</sup>

Erasmus as an author and a consultant worked together with the Venetian printer and publisher several times. He already wrote to Aldus, that his translations of Euripides will make him immortal, especially if they will be "printed using your small [minutius] types, which are the most elegant in the world (*tuis excusae formulis... maxima minutioribus illis omnium nitidissimis*)."<sup>7</sup> But beauty is not everything. Erasmus calls the octavo editions the outstanding products of his age. Not even Ptolemaios of Philadelphia could access literature and science the way Aldus Manutius made it possible in the early modern age. While the great king had built only one extensive library, Aldus raised 'a library without walls,' which will survive all disasters. Willibald Pirckheimer, yet another humanist, esteemed his Theokritos-edition so highly, that he asked the famous German engraver and painter, Albrecht Dürer to illustrate its cover and design an *ex libris* sign for it. Every detail of the bucolic idyll created by Dürer follows the text faithfully.

There was feedback from Hungary, too. Sigismund Thurzo provost and secretary of the king writes in his grateful letter of 1501: "My different kinds of affaires con-

<sup>5</sup> Stefan Zweig: *The Glory and Tragedy of Erasmus of Rotterdam*. Budapest, 1993. Holnap Kiadó. pp. 73–80.

<sup>6</sup> Nándor Várkonyi: *The History of Books and Letters*. Budapest, 2001, Szépirodalom Könyvműhely. p. 352.

<sup>7</sup> Anthony Grafton: *Humanist Reading*. In: *Cultural History of Reading*. Budapest, 2000, Balassi. Eds. Robert Chartier, Guglielmo Cavallo. p. 206.

sume the time I could spend at home in the company of poets and orators. Your books – being very practical so that I can take them with me for my walks or have them around during conversations or my affaires in the court – cause me much delight.”<sup>7</sup>

No doubt, Aldus Manutius was the leading publisher in Europe at the turn of the 15<sup>th</sup> and 16<sup>th</sup> centuries. A generation after Gutenberg, printers overstepped the traditional manuscripts of medieval times and by appeasing the needs of this new generation of readers, published Latin and Greek classics with flawless content and in wonderful form.

Aldus died in 1515, but his workshop continued to operate throughout the 16<sup>th</sup> century. His work, if not to the same effect but still at a very good standard, was carried on by his son and later by his grandchild for a hundred years. In 1597 the grandchild gave up the workshop in Venice and in response to the call of the pope he went to Rome to manage the printing facilities of the Vatican. By that time the humanistic movement was a thing of the past. Public opinion was mainly concerned by the struggle of the catholic church with protestants, and its efforts to renew itself.<sup>8</sup>

As they were popular the number of surviving Aldine books is not too great and oftentimes they are worn-out because of frequent usage. But their quality is clearly shown by the fact that they were still good enough for the famous French playwright Racine (1639–1699) who got to know the Greek tragedies from Aldine editions nearly a hundred and fifty years after their publication.

Finally let us attempt to answer the question in the title of this paper. As we have seen Aldus was seriously into the business of publishing ‘scientific’ books. He was open to the

novelties of his trade, moreover he also made significant contributions to it with his inventions. Based on this we could assume that if computers, the main representatives of modernity and future had existed at his time, he would have surely used them. But would he have used T<sub>E</sub>X?

Aldus published many works in Greek. (By the way the word T<sub>E</sub>X comes from the Greek ΤΕΧΝΗ or techné). Originally T<sub>E</sub>X being an American software had a limit of working with 128 characters only. This could have raised initial problems for someone planning to bring out multi-lingual publications. However, this limit has been eliminated since then, and at this conference we hear about the both typographically and technologically difficult task of publishing the Koran itself.

Though I am not a T<sub>E</sub>X guru or a T<sub>E</sub>Xpert, I understand the many advantages of T<sub>E</sub>X in the fields of scientific publications especially if it contains a lot of formulae. Well, we must realize that the early modern age of Aldus preceded Newton’s scientific revolution. In the 16<sup>th</sup> century we could have possibly find some formula in esoteric works by alchemists on how to create gold, but usually explained in very unclear terms. As far as the ‘serious science’ of the century is concerned, it was mainly represented by classical works of ancient authors. They did not perform experiments back then but read Aristotle instead, so science was more like literature.

We can state that Aldus’ innovative personality was always ready to accept newer and better solutions. He would have thought about using T<sub>E</sub>X if Donald Knuth had been born several hundred years earlier. But since the works he published were not scientific but rather literary in form, he might have decided otherwise.

<sup>8</sup> Lucien Febvre and Henri-Jean Martin: *The Coming of the Book (L’Apparition du livre). The Impact of Printing 1450–1800*. London, 1990, Verso. p. 124.





# Practical T<sub>E</sub>X 2006:

## L<sup>A</sup>T<sub>E</sub>X Workshop and Presentations on

### L<sup>A</sup>T<sub>E</sub>X, T<sub>E</sub>X, ConT<sub>E</sub>Xt, and more

L<sup>A</sup>T<sub>E</sub>X workshop: July 25–28, 2006

Practical T<sub>E</sub>X conference: July 30–August 1, 2006

Rutgers, the State University (Busch Campus)  
Piscataway, New Jersey, USA

<http://tug.org/practicaltex2006>  
[conferences@tug.org](mailto:conferences@tug.org)

**Keynote address:** *Barbara Beeton*,  
American Mathematical Society & T<sub>E</sub>X Users Group

- April 1, 2006 - presentation proposal deadline
- April 15, 2006 - early bird registration deadline
- July 14, 2006 - hotel reservation deadline

#### Further information

Conference attendees will enjoy an opening night reception and an (optional) banquet on the final evening. Coffee and lunch will be served each day of the meeting. Located on the Busch Campus of Rutgers University in Piscataway, New Jersey, an easy train ride from New York City.

Conference fee, hotel, and other information is available on the web site.

#### Pre-conference L<sup>A</sup>T<sub>E</sub>X workshop

A dedicated hands-on workshop covering a significant portion of L<sup>A</sup>T<sub>E</sub>X will be offered for the four days before the conference (Tuesday–Friday). Time permitting, the instructors will help with individual projects. Class space is strictly limited, so please reserve early.

#### Sponsorship

If you'd like to support the conference, promote your T<sub>E</sub>X products and services, or otherwise provide sponsorship, see the web site for donation, sponsorship, and advertising options.

We are very grateful to Rutgers University for major support, the German-speaking T<sub>E</sub>X users group DANTE for a special contribution, and especially to the many individual contributors.

Hope to see you there!

*Sponsored by the T<sub>E</sub>X Users Group.*





# TUG 2006

## Digital Typography & Electronic Publishing: Localization & Internationalization

*The 27<sup>th</sup> Annual Meeting & International Conference of the T<sub>E</sub>X Users Group*  
Announcement and Call for Papers



After TUG 2003 in America (Hawaii, USA), TUG 2004 in Europe (Xanthi, Greece), TUG 2005 in Asia (Wuhan, China), it is fitting that TUG 2006 be held in Africa—in Marrakesh, Morocco. The conference will be hosted by Cadi Ayyad University's Sciences Faculty.

Processing multilingual e-documents will go beyond the limits of its traditional cultural areas and new horizons in the internationalization of T<sub>E</sub>X will be explored. The conference consists of scientific and technical talks, panels and poster sessions, as well as tutorials, all dedicated to present investigations of e-document technology.

The topics of interest include:

### *System design:*

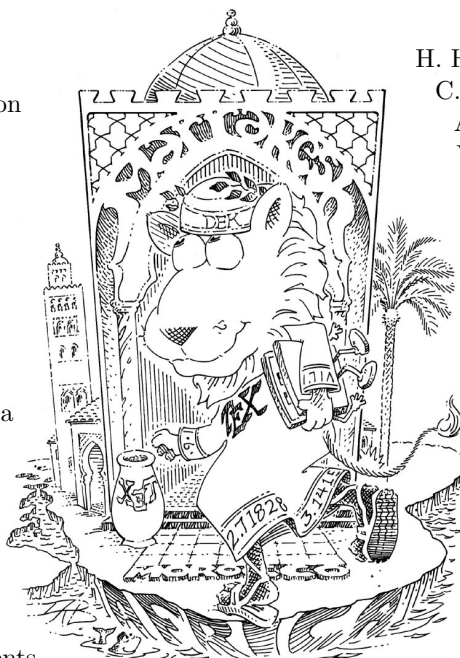
Internationalization & localization  
Presentation and content  
Electronic publishing  
Digital typography  
Document design  
Recognition  
Calligraphy  
Standards

### *Web technology:*

Web semantics & structured data  
Multimedia and hypertext  
Structure transformations  
Digital libraries

### *Educational technology:*

Interactive e-learning  
Learning support systems  
Learning and dynamic e-documents



### *Program committee:*

H. Hagen, Pragma ADE, Netherlands  
C. Swanepoel, UNISA, South Africa  
A. Lazrek, UCAM-FSSM, Morocco  
V.R.W. Schaa, DANTE, Germany  
Y. Haralambous, ENSTB, France  
K. Sami, UCAM-FSSM, Morocco  
S. Peter, Beech Stave Press, USA  
K. Höppner, DANTE, Germany  
A. Syropoulos, GTF, Greece  
J. Plaice, UNSW, Australia  
A. Lindsay, Lincs Uni, UK  
B. Hughes, UM, Australia  
B. Raichle, US, Germany  
C. Beccari, PT, Italy  
K. Berry, TUG, USA

### *Organizing committee:*

K. Berry, TUG, USA  
M. El Adnani – A. Lazrek –  
K. Sami, UCAM-FSSM, Morocco

Tutorial days: *November 7–8, 2006*

*July 30, 2006:* Paper submission

Main conference: *November 9–11, 2006*

*August 30, 2006:* Early registration

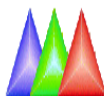
For more information, paper submission, and online registration

please visit: <http://www.tug.org/tug2006>

or email: [tug2006@tug.org](mailto:tug2006@tug.org)



UCAM



STIC



GUTenberg