

# The Components of T<sub>E</sub>X

Joachim Schrod  
Detig · Schrod T<sub>E</sub>Xsys

March 1991

## Abstract

T<sub>E</sub>X needs a great amount of supplementary components (files and programs) of which the meaning and interaction often is unknown. This paper explains the components of the kernel system T<sub>E</sub>X that are visible for the T<sub>E</sub>X user and their relations.

## 1 About this Report

T<sub>E</sub>X is a typesetting system which offers authors easy usage of powerful typesetting features to produce printed matter which is the state of the art of computer typesetting. This is, however, not done by the T<sub>E</sub>X program alone: A significant number of supplementary programs and files together form the complete typesetting and authoring system. Along with the programs that belong to T<sub>E</sub>X directly, there exist two other major programs which were built by DONALD KNUTH in connection with T<sub>E</sub>X and must be included in an explanation of the full system: METAFONT, for the generation of fonts, and WEB, a documentation and developing ‘language’ for programming. T<sub>E</sub>X and METAFONT are written in WEB.

This text describes this ‘kernel’ T<sub>E</sub>X from a user’s viewpoint: at the end you should have an overview of the ingredients of the T<sub>E</sub>X system, and about the files and support programs that are essential for you as a user. This will not, however, be an introduction to the capabilities of T<sub>E</sub>X or how you may run T<sub>E</sub>X on your computer.

I will use marginal notes to identify the places where terms are explained for the first time. Abbreviations for file types – usually identified by common suffixes or extensions – are set in a monospace type (“typewriter”) and those abbreviations are put into the margin, too. Please note that these abbreviations are sometimes not identical with the file extensions (see also Table 1).

This report is the start of a series that describes the

subsystems mentioned above and their respective components. In that series each report will focus on one subsystem in a special point of view; it should not result in gigantic descriptions which tell everything (and then nothing). In my opinion the following reports will be of interest:

- the structure of a standard installation of T<sub>E</sub>X
- DVI drivers and fonts
- possibilities of graphics inclusion in T<sub>E</sub>X documents
- the components of METAFONT
- the structure of a standard installation of METAFONT
- WEB systems — the concept of Literate Programming
- other (though not yet planned) themes of interest are perhaps
  - differences between T<sub>E</sub>X and DTP systems
  - the way how T<sub>E</sub>X works (there exists some good books on this topic!)
  - the limits of T<sub>E</sub>X
  - T<sub>E</sub>X as a programming language

The reports will be published in this sequence.

## 2 What is T<sub>E</sub>X?

T<sub>E</sub>X is a typesetting system with great power for the typesetting of formulae. Its basic principle is that structures in the document are marked and transformed into typeset output. Providing such information about the structure of a document is known as *markup*. If the marks describe the look of the document, it is called *optical markup*, while, if document structures are marked, it is called *logical markup*. T<sub>E</sub>X provides both forms of

---

© Copyright 1990, 1991 by Joachim Schrod. All rights reserved. Republication and Distribution is allowed under conditions analog to the GNU General Public License. Previous revisions of this paper appeared in “Die T<sub>E</sub>Xnische Komödie” and in “Baskerville.”

Actual address: Joachim Schrod, Detig · Schrod T<sub>E</sub>Xsys, Kranichweg 1, D-6074 Rödermark, FR Germany, Email: xitijsch@ddathd21.bitnet

markup, i. e., exact control of the layout of parts of the document and their positioning as well as the markup of the structure of formulae or document components. The logical markup is mapped to the optical one by T<sub>E</sub>X so that layout may serve for the identification of structures by the reader.<sup>1</sup>

The kernel of the T<sub>E</sub>X typesetting system is the formatting program T<sub>E</sub>X82, which is often simply called T<sub>E</sub>X. This usage shall be adopted here whenever the difference between the complete system and the formatter is unimportant or obvious. T<sub>E</sub>X82 is a big monolithic program which is published in the book *T<sub>E</sub>X: The Program* by DONALD KNUTH. Its features may be separated into two levels:

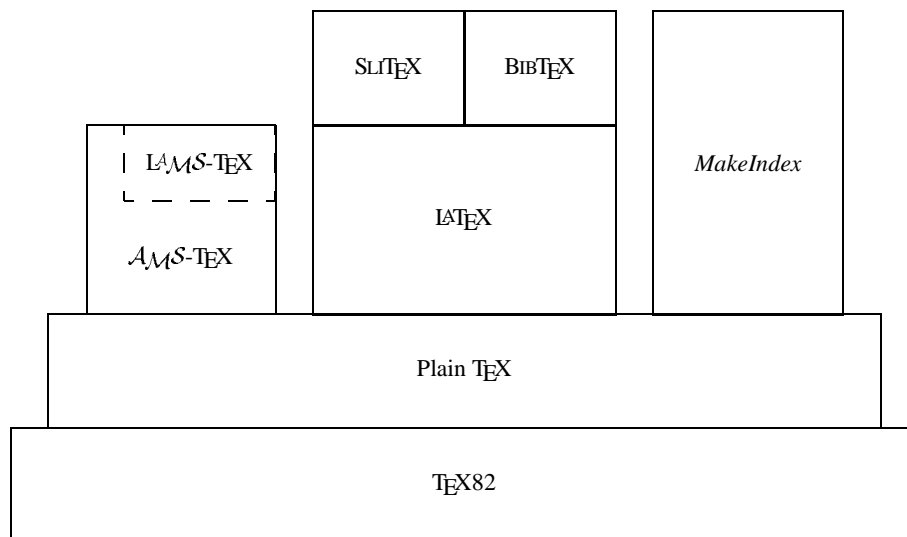
1. T<sub>E</sub>X82 formats text, i. e., it breaks it into paragraphs (including automatic hyphenation) and produces page breaks.
2. It provides the programming language T<sub>E</sub>X which incorporates a macro mechanism. This allows new commands to be built to support markup at a higher level. DONALD KNUTH presents an example in the T<sub>E</sub>Xbook: Plain T<sub>E</sub>X. A collection of macros which supports a special task and has (hopefully) a common philosophy of usage is called a *macro package*.

High level features for optical markup, as represented by Plain T<sub>E</sub>X, allow one to build additional levels leading to full logical markup. At the moment, two

macro packages for logical markup are widespread:  $\mathcal{A}\mathcal{M}\mathcal{S}$ -T<sub>E</sub>X and L<sup>A</sup>T<sub>E</sub>X. Both systems are built on top of Plain T<sub>E</sub>X to greater or lesser extents and the user can use the optical markup of Plain T<sub>E</sub>X in addition to logical markup if desired. This results in the effect that the author can use a mixture of structural information and explicit layout information – a situation with a high potency of features that nevertheless can (and does) lead to a lot of typographic nonsense!

As T<sub>E</sub>X82 was built only for typesetting texts and to allow the realization of new markup structures, many features are lacking which are required by authors. To provide features like the production of an index or a bibliography or the inclusion of graphics, additional programs have been written, which use information from a T<sub>E</sub>X82 formatting run, process them, and provide them for the next T<sub>E</sub>X82 run. Two supplementary programs are in widespread use and available for many computer/operating system combinations: BIBT<sub>E</sub>X, for the production of a bibliography from a reference collection, and *MakeIndex*, for the production of an index.

A special case of the processing of information provided by a T<sub>E</sub>X run is the production of a table of contents or the usage of cross references in a text. For this only informations about page numbers, section numbers, etc., are needed. These are provided by T<sub>E</sub>X82 and can be processed by T<sub>E</sub>X82 itself, so T<sub>E</sub>X82 is used as its own post processor in this situation.



Figuur 1: The Components of T<sub>E</sub>X

<sup>1</sup>Layout – and book design in general – do not represent useless beauty. A good book design must first support the understanding of the content to produce readable text. So it is *aesthetic* in its best sense, since it connects form and contents and builds a new quality.

We have now seen that the T<sub>E</sub>X typesetting system is a collection of tools that consists of the typesetting ‘engine’ T<sub>E</sub>X82, macro packages (maybe several that are based on others) and supplementary programs, used together with these macro packages. This relation is illustrated by Figure 1.

### 3 Formatting

The formatting process of T<sub>E</sub>X needs information about the dimensions of characters used for the paragraph breaking. A set of characters is grouped in *fonts*. (But this is a simplification as the notion “font” should be used for the realization of a type in a fixed size for a specific output device.) The dimensions of the characters of a font are called *font metrics*.

The format in which the font metrics are used by T<sub>E</sub>X was defined by DONALD KNUTH and is called TFM format (“T<sub>E</sub>X font metrics”). In this format, every character is described as a *box* with a height, a depth, and a width. T<sub>E</sub>X only needs these measurements, it is not interested in the shape of the character. It is even possible that the character may extend outside the box, which may result in an overlap with other characters. The character measures are specified in a device independent dimension because T<sub>E</sub>X processes its breaking algorithm independent of any output device.

During paragraph breaking, T<sub>E</sub>X hyphenates automatically, which can be done in an almost language-independent way. For the adaption to different languages, *hyphenation patterns* are needed to parametrize the hyphenation algorithm.

The result of a T<sub>E</sub>X formatting run is a DVI document, in which the type and position on the page are specified for each character to be output. The resolution that is used is so small that every possible output device will have a coarser raster, so that the positioning is effectively device independent. The DVI document specifies only types, not the fonts themselves, so that the name DVI<sup>2</sup> (“*device independent*”) is accurate. To make the result of the formatting run available, the DVI file must be output by a so-called DVI driver on the desired output device.

If problems occur during the formatting, error messages or warnings are output on the terminal. *Every* message that appears on the terminal will also be written into a protocol file named LOG file. In this LOG file additional information may be placed that would have been too verbose for the output to the terminal. If this is the case, T<sub>E</sub>X will tell the user so at the end of the formatting run. The messages of T<sub>E</sub>X are not built in the program, they are stored in a (string) POOL file. These messages must be read in at the beginning of a run.

<sup>2</sup>This name is a problem because “DVI” is a trademark of Intel Corp. now, but the name DVI for T<sub>E</sub>X output files pre-dates this.

### 4 Macro Packages

The basic macro package is Plain T<sub>E</sub>X, developed by DONALD KNUTH together with T<sub>E</sub>X82. It parametrizes the T<sub>E</sub>X82 typesetting machine so that it can typeset English texts with the Computer Modern type family. Additionally, Plain T<sub>E</sub>X provides optical markup features. Plain T<sub>E</sub>X is available as one source file, `plain.tex`.

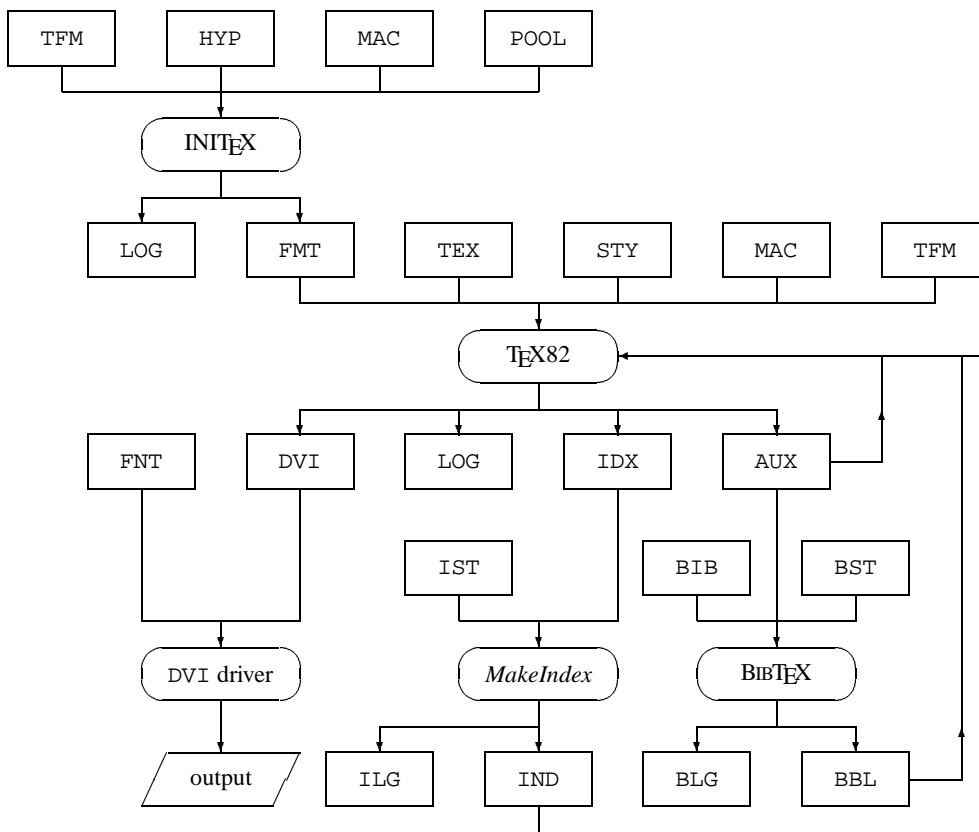
All other macro packages known to the author are based on Plain T<sub>E</sub>X, i. e., they contain the source file `plain.tex` either originally or with modifications of less important parts. Next to Plain T<sub>E</sub>X, the most important (free) available macro packages are  $\mathcal{A}\mathcal{M}\mathcal{S}$ -T<sub>E</sub>X by MICHAEL SPIVAK and L<sup>A</sup>T<sub>E</sub>X by LESLIE LAMPORT. Other free macro packages are often of only local importance (e. g. BlueT<sub>E</sub>X, TEXT1, or T<sub>E</sub>Xsis) or are used in very special environments only (e. g. `texinfo` in the GNU project or `webmac` for WEB). Important commercial macro packages are MacroT<sub>E</sub>X by AMY HENDRICKSON and L<sup>A</sup>M<sup>S</sup>-T<sub>E</sub>X, also written by MICHAEL SPIVAK.

These macro packages usually consist of a kernel that provides additional markup primitives. With such primitives, *document styles* can be built which realize logical markups by a corresponding layout. This layout can often be varied by *sub-styles* or *style options* which may also provide additional markups.

The macro packages produce *supplementary files* which contain information about the page breaks or the document markup. This information may be used by support programs – e. g., the specification of a reference from a bibliography database or the specification of an index entry with corresponding page number for the construction of an index. A special case is the information about cross references and headings for the building of a table of contents, as this information can be gathered and re-used by T<sub>E</sub>X directly.

SL<sub>T</sub>E<sub>X</sub> is a special component of L<sup>A</sup>T<sub>E</sub>X for the preparation of slides with overlays. In TUGBoat volume 10, no. 3 (1989) L<sup>A</sup>M<sup>S</sup>-T<sub>E</sub>X was announced, which will provide the functionality of L<sup>A</sup>T<sub>E</sub>X within  $\mathcal{A}\mathcal{M}\mathcal{S}$ -T<sub>E</sub>X. MacroT<sub>E</sub>X is a toolbox of macro “modules” which may be used to realize new markups but, as it became available only short time ago, it is not yet widespread.

For the usage of these (and other) macro packages, one must check whether they need additional fonts which do not belong to the Computer Modern type family. For L<sup>A</sup>T<sub>E</sub>X, e. g., fonts with additional symbols and with invisible characters (for the slide overlays) are needed, while  $\mathcal{A}\mathcal{M}\mathcal{S}$ -T<sub>E</sub>X needs several additional font sets with mathematical and Cyrillic characters.



Figuur 2: The Connection of Components and File Types

## 5 Support Programs

Only two support programs will be discussed here: **BIBTEX** by OREN PATASHNIK for the preparation of bibliographies and **MakeIndex** by PEHONG CHEN and MICHAEL HARRISON for the preparation of a sorted index. For both tasks exist other, functionally equivalent, support programs. But the abovementioned are available on many operating systems, and have an “official” state as they are LESLIE LAMPORTE encourages their usage with **L<sup>A</sup>T<sub>E</sub>X** in his documentation, and the TUG supports them for general use.

There is no totally portable mechanism for the inclusion of general graphics in T<sub>E</sub>X documents, so that there are no machine independent support programs available.

**BIBTEX** is used to handle references collected in **BIB** files. T<sub>E</sub>X produces supplementary files which contain information about the required references, and **BIBTEX** generates from them a sorted bibliography in a **BBL** file which may be subsequently used by T<sub>E</sub>X. The kind of sorting and the type of cite keys are defined by *bibliography styles*, specified in **BST** files. The messages of a **BIBTEX** run are written to a **BLG** logfile.

**MakeIndex** reads an **IDX** support file that contains the index entries and the according page numbers, sorts

these items, unifies them and writes them as T<sub>E</sub>X input in an **IND** file. The formatting style may be specified by an *index style*. The messages of a **MakeIndex** run are written to a **ILG** file.

## 6 Performance Improvements

Much of the work that T<sub>E</sub>X82 has to do is the same for every document:

1. All text has to be broken into lines. Text pieces in the same language are hyphenated with the same hyphenation patterns.
2. The basic markups of the corresponding macro packages must be available.
3. The required font metrics are much alike for many documents, as the font set used usually doesn't differ that much.

To improve T<sub>E</sub>X's performance, hyphenation, markup, and font metrics descriptions are converted from an external, for (1) and (2) textual, representation into an internal representation which can easily be used by T<sub>E</sub>X82. It is sensible to do this transformation only once, not for every document. The internal representation is stored in a **FMT** file. The storing is done with the T<sub>E</sub>X command `\dump`, so that **FMT** files often are

FILE TYPE	EXPLANATION	FILE IDENTIFICATION (SUFFIX, EXTENSION, ETC.)
TEX	Text input	tex, ltx
DVI	$\TeX$ 82 output, formatted text	dvi
LOG	$\TeX$ 82 log file	log, lis, list
HYP	Hyphenation patterns	tex
TFM	Font metrics	tfm
POOL	String pool	pool, poo, pol
FMT	Format file	fmt
MAC	$\TeX$ macro file	tex, doc
STY	$\TeX$ style file	sty, tex, st, doc
AUX	Support files	aux, toc, lot, lof, glo, tmp, tex
BIB	Reference collections	bib
BBL	References or bibliographies	bbl
BLG	BIB $\TeX$ log file	blg
BST	BIB $\TeX$ style file	bst
IDX	Unsorted index	idx
IND	Sorted index	ind
IST	Index markup specification	
ILG	<i>MakeIndex</i> log file	ilg

Tabel 1: *File Types*

called “*dumped formats*.” A FMT file can be read at the beginning of a  $\TeX$ 82 run and is thus available for the processing of the actual text.

As the creation of a FMT file is done infrequently – usually for the update of a macro package – the formatting of texts can be done with a reduced version of the  $\TeX$ 82 program that doesn’t contain the storage and the program parts for the transformation of the hyphenation patterns and for the dumping. The complete version of  $\TeX$ 82 is needed in an initialization phase only and therefore called INIT $\TeX$ . Additional improvements of the performance can be reached by the usage of production versions of  $\TeX$ 82 from which parts for statistical analysis and for debugging are stripped.

$\TeX$  versions that have no dumped formats preloaded, have the ability to load a dumped format (i. e. a FMT file), and have no ability to dump a FMT file (i. e., they are not INIT $\TeX$ ) are often called Vir $\TeX$ , which stands for *virgin  $\TeX$* .

## 7 Connections Between File Types and Components

In the above sections, the components of the  $\TeX$  authoring system were described, and the files that are read or written by these components mentioned. The connections between them all is demonstrated graphically in Figure 2. In this graphic, file types are represented by rectangles, and programs by ovals. The arrows mean “is read by” or “is produced by.” The abbreviations of the file types are explained in Table 1, which also lists the file identifications (suffixes or extensions) that these files usually have (but note that other file identifications are also in use).

### Acknowledgements

I would like to thank CHRISTINE DETIG who was so kind to provide the English translation. NELSON BEEBE suggested performing the translation. KLAUS GUNTERMANN made valuable comments on the first (German) version. NICO POPPELIER contributed a new version of figure 2, better than my original one.