

A brief history of T_EX, volume II

Arthur Reutenauer*

ENST Bretagne

Technopôle Brest-Iroise

CS 83818

29238 BREST CEDEX 3

arthur.reutenauer@normalesup.org

0 Rationale for this *volume II*

When I gave this talk in Bachotek I appended the subtitle *Pax T_EXnica—the program on which the sun never sets*—an obvious pun to two historical empires renowned for their considerable geographical extent. I wanted to add this subtitle for two reasons: first, I liked *A brief history of T_EX* a lot but I realized after choosing it that there had already been a talk with this exact same title more than ten years before² and I wanted to avoid the risk of confusion, be it only for archival purposes; and second, I felt the subtitle made my standpoint clear: the history I wanted to account for was very much a *geographical* one: how T_EX enabled us to gradually typeset in every language of the world—or almost so. As far as the printed version was concerned, though, it seemed that it could also be considered a sequel to the first article—after all, many things had changed over ten years! Hence this *volume II*.

But first, let us recapitulate things from the beginning . . .

1 The origins

1.1 In the beginning there was . . .

History begins, scholars tell us, with the invention of writing and the ability to account for one’s own culture. So, in the beginning there was typesetting and the program that enabled us to do so, let us call it T_EX. This program was written by a man, let us call him 高德纳. Oh, and we need a date, too, so let’s say 1978, thirty years ago.

高德纳, as its name suggests, lived in a region inhabited by many Chinese citizens, next to the great city that goes by the name of the Old Golden Hills. But he was an American citizen and a native speaker of the English language. So the program he

* I wish to thank Jerzy Ludwiczowski heartily for his constant encouragements to write this article and his patience in waiting for it.

² The talk had been held in Toruń in 1995 and published the next year by different journals, including the *TUGboat* where it is available online: <http://www.tug.org/TUGboat/Articles/tb17-4/tb53tay1.pdf>.

Figure 1: The name of a distant galaxy. From the very beginning, T_EX sets out to conquer the universe (extract from `story.tex`, in *The T_EXbook*, chapter 6).

wrote was all in English (with a lot of ‘\’ though) and it was meant—at first—for English speakers to use.

Nevertheless, when 高德纳 created T_EX, he still thought of the users speaking other languages. Of course, all the commands were in English, the default settings were chosen for that language, and the fonts used a 7-bit encoding³ supporting only the Latin script, but he made provisions for extending this. The fonts, in particular, were supplied with a set of diacritics with the help of which he devised the well-known accent commands that could construct an accented character “on the fly”, a sample of which can be seen on figure 1, an extract of a famous T_EX file. This enabled to write, mostly, all the languages of Western Europe—and therefore of all the other countries that use the same languages, in particular the entire South America.

This was the first step since even if they may seem unpractical, the accent commands actually introduced a way of inputting a lot of characters users didn’t have access to on a standard American keyboard, much in the way math commands were a way of specifying the layout of complicated math formulæ; so even if they were not an encoding⁴ in the current meaning of the term, they were a sort of coding system, and were thought as such by many users as well as some recoding utilities⁵.

³ That is, they could only have up to $2^7 = 128$ characters.

⁴ In the sense of an encoded character set, like ASCII, the ISO-8859-* family, or Unicode.

⁵ To name but only one, the very popular `recode` program (<ftp://ftp.gnu.org/pub/gnu/recode/>) knows about sequences like `\’e` as “T_EX” encoding.

“Uznając, iż los nas wszystkich od ugruntowania i wydoskonalenia konstytucji narodowej jedynie zawisł, długim doświadczeniem poznawszy zadawnione rządu naszego wady, a chcąc korzystać z pory, w jakiej się Europa znajduje i z tej dogorywającej chwili . . .”

Figure 2: Polish uses a lot of diacritics (extract from the May Third Constitution, 1791)

So \TeX extended, from the very beginning, over all the Americas as well as the western part of Europe, and many regions in Africa⁶.

But, as mentioned, this wasn’t enough even for some other languages using the Latin alphabet, let alone languages using any other alphabet or different writing systems. Work needed to be done, as 高德纳 acknowledged that he couldn’t handle all the languages of the world by himself, and he encouraged people to settle to this task. It wasn’t long before people did indeed so.

1.2 Go East

As \TeX was born, the story tells us further, there was a companion program called METAFONT, whose purpose was to design the fonts that \TeX used. As a matter of fact, all the letters and accents we discussed above were all drawn using METAFONT, so adapting them the fonts to other languages meant, mostly, drawing more characters as needed.

Let’s start how this was done. An interesting example is Polish. It uses, too, a wealth of accents (see figure 2); most of them were already present in the fonts or easy to add, by simple modifications to the existing characters. One of these accents, though, is quite special: it’s called *ogonek* which means “little tail” in Polish, as for example on the first word of figure 2 It looks remotely like a symmetrised cedilla but not quite, and the drawing had therefore to be invented entirely and polished carefully⁷. Then, a new control sequence had to be invented and agreed upon in order for users to be able to input that characters, in the same spirit as the already existing accent command; nowadays, it’s $\backslash\mathbf{k}$ in \LaTeX ⁸.

⁶ Including, of course, all the languages of the former colonizers like English and French, but also important African languages like Swahili which are written entirely in the Latin alphabet.

⁷ The Poles are very proud of their *ogonek* and you should not upset them by speaking ill of it. Maybe it is even too daring in the eyes of some to state that *ogonek* looks like an inverted cedilla!

⁸ For a thrilling account of how \TeX came to Poland, I heartily recommend to read the text of this talk, given at the TUG meeting in Hawai’i in 2003, and published in the *TUGboat*, volume 24, No. 1: <http://www.tug.org/TUGboat/Articles/tb24-1/odyniec.pdf>.

Over the years, more characters were designed and entire alphabets were digitized using METAFONT, starting with Greek and Cyrillic, which were drawn by various people around the world.

An important step was when \TeX was extended in 1989 to handle 8-bit (then becoming $\TeX3$), thus enabling fonts to have up to 256 characters. The next year, during a meeting in Cork, \TeX users from all over the world agreed on a standard encoding for \TeX ’ latin fonts, which then came to bear the name of its birth place (or the alternative, less poetical names of $\mathbf{T1}$ or $\mathbf{8t}$). Another important stepmark at that time was the advent of the \LaTeX *babel* package, which attempted to provide a convenient way to switch between languages and a common interface for the \LaTeX users.

But even after those fonts were designed, after those standards were agreed upon, many things were left to do: what about Arabic, for example? \TeX offered amazing possibilities, but did not really address the issues of right-of-left typesetting and it also completely left aside the fact that characters can have different forms according to their place in a word (both being essential features of Arabic). Therefore, to go further it was necessary to *think different*⁹!

2 Think different

2.1 \TeX encompasses the *Mare Nostrum*

As soon as 1987, the first experiments were made to handle the challenge of Arabic typesetting and gave birth to a modified version of \TeX called $\TeX\text{-X}\mathbf{\E}\mathbf{T}$, to emphasize the fact that it could write in two different directions¹⁰. This worked in a particular way: when writing data in the output file, $\TeX\text{-X}\mathbf{\E}\mathbf{T}$ did not reverse the order of letters but put a mark whenever it encountered a sequence of Arabic letters, and let all the work be done by the printer driver. That way, the text was stored in *natural* order in the output file—that is in the order in which an Arabic speaker would speak out the letters—but, on the other hand, it meant the files output by $\TeX\text{-X}\mathbf{\E}\mathbf{T}$ had to be processed by a special driver. So 高德纳—who, once again, was the

⁹ To quote an old slogan of one of the big computer manufacturers—I’m not sure what the legal status of such commercial slogans is and I may not be entitled to reuse it in a document; but I want to make sure people know I didn’t mean any harm and in case TUG is sued I deny everything.

¹⁰ The founding article was published in the *TUGboat*, volume 8, No. 1: <http://www.tug.org/TUGboat/Articles/tb08-1/tb17knutmix.pdf> and makes a fascinating reading even today, especially when compared with the current paradigm established by Unicode in that area—the so-called bidirectional algorithm.

input: Book is "باتك" in Arabic
output: Book is "كتاب" in Arabic

Figure 3: The challenge of Arabic writing: when setting an Arabic text, the order of the letters does not only need to be reversed, but their shapes also may vary a lot—can you recognize all four of them on the second line?

lead in that project—decided that the output format should be called DVI-IVD to differentiate it from the traditional DVI output format.

That way, both the \TeX program and the DVI were “extended” in the sense that they were made able to handle different types of information in addition to the ones they already knew how to process or store. We shall meet a lot of these along the way, and we shall refer to them as “extensions”—or sometimes “engines” for \TeX extensions. So $\text{\TeX-X}\epsilon\Gamma$ was, probably, the very first \TeX extension.

A few years later, $\text{\TeX-X}\epsilon\Gamma$ was itself extended in something that achieved roughly the same goals, but without needing to resort to an extension of the DVI: it readily reversed the order of each letter in the output file as appropriate. To mark both the similarities and the difference of this second extension with the first one, it was called by the same but a second hyphen: $\text{\TeX--X}\epsilon\Gamma$!

These improvements were interesting and made Arabic typesetting with \TeX possible very early; but it was still an experimental stage, and apart from that, it did not change things for other scripts like, in particular, the Indic and South-Asian scripts.

2.2 Enters Unicode

For a better suited treatment of such complex scripts, Ω ¹¹ was designed. It consisted of several major improvements:

- It enabled (probably) every sort of writing directions.
- It came with a set of filters (the Ω transformation processes, ΩTP for short) that transformed the input text.
- It enhanced the traditional font formats used by \TeX from 8-bit based encoding to 16 bits.

The two first points made the treatment of Arabic much more natural (just switch the writing direction from left-to-right to right-to-left, top-to-bottom; and filter the input text to give each letter its appropriate appearance given the context); and the third one was also a very important one because it addressed the problem which we haven’t yet men-

¹¹ We shall call it simply by the Greek letter from now on.

tioned: up to then, \TeX handled only fonts with at most 256 slots, which wasn’t so important for alphabetic scripts but became a major issue when one wanted to typeset in a language using ideographs, whose number by far exceeds this number.

Ω lifted a part of this problem by making direct use of (possibly) very large font metrics; that is it could use any font on the input but remained constrained by the output format.

Anyway, it brought with it a conceptual leap, even if it failed to address some of the issues of the output format. Over the years it has been successfully used to typeset the Devanāgarī, Malayalam, Tibetan, Inuktitut and Cherokee scripts among others, although it has never really gained a wide acceptance.

2.3 The Other Way: generating PDF

Ω was first formally released in 1994, and by that time there was a document format that was increasingly gaining in popularity and commercial strength: PDF. Hearing about this “Portable Document Format” in the \TeX world, one cannot help thinking that it is a concept quite close to the traditional output format, DVI (does it not stand for “DeVice-Independent?”); therefore it seemed only right that \TeX should be able to produce PDF directly: and so it did, with the birth of the well-known pdf \TeX on March 15th, 1997 (then under the name `tex2pdf`).

Another huge improvement pdf \TeX brought was the direct handling of TrueType fonts, which had that time had become the major font format for personal computers.

2.4 One more extension ...

Worth mentioning here, since its later development was to be closely related to pdf \TeX ’s, is “the” extension of \TeX , called $\epsilon\text{-}\text{\TeX}$ for “extended \TeX ”. Developed during the late 90ies, it extended $\text{\TeX-X}\epsilon\Gamma$ which we mentioned earlier (the second one, with two hyphens) and was therefore of great use of Arabists and other communities writing from right to left.

Its very useful extended features were later merged into pdf \TeX which had for a while an offspring called pdf $\epsilon\text{-}\text{\TeX}$, now fully incorporated into pdf \TeX ; that is, pdf \TeX now supports the $\epsilon\text{-}\text{\TeX}$ extensions, but it can also pretend to know nothing about these and be simply pdf \TeX ¹².

¹² Just as it has also could, from the very beginning, behave as the DVI-producing \TeX , or the actual pdf \TeX —which means that in DVI mode and *with* the extensions ... you probably get the picture.

2.5 Needless To Say

Before proceeding to the last part of this account, there are a few words to be said on another attempt of extending \TeX , which isn't very famous now but whose name still lingers in many memories. It was to be a completely new concept, opening up a world of possibilities ... as of today, it is but no more.

The “New Typesetting System”, so it was called—or $\mathcal{N}\mathcal{T}\mathcal{S}$ for short—was a complete reimplementa-tion of \TeX in Java, aiming at full compatibility with the original engine, and providing at the same time the great modularity and extensibility that comes with that language.

Sadly, while the first goal was actually achieved (\TeX was indeed rewritten in Java), it proved completely unusable and pointless because of its extreme slowness. The extension projects were never carried through and $\mathcal{N}\mathcal{T}\mathcal{S}$ has now been officially declared dead.

3 Rule, \TeX annia, \TeX annia, rule the waves

3.1 Taming the multilingual lion

Back to living projects now, there is one very young which has brought many changes for a lot of users recently. called, is a very multilingual extension of \TeX ; the very name suggests, again, that it can type in every direction (it can spell “ \TeX ” backward). Its spirit is a bit special in that it started off (in April 2004) as a Mac OS-specific programm which made heavy use of the Apple libraries designed to handle text and scripts¹³.

But shortly after (April 29th, 2006), $\text{X}\mathcal{F}\text{\TeX}$ was released for Linux too, and it was not long before it was ported to Windows as well.

The huge difference of $\text{X}\mathcal{F}\text{\TeX}$ —and overwhelming advantage for many of its adepts—was to get rid of practically all the hassle in font selection, font installation, etc., while opening up at the same time a whole new world of possibilities: people could at once use all the bleeding-edge features of the newest font technologies with no particular problem. The key to this was the use of a lot of external libraries—which of course comes at a price: users lost part of the control they had on every detail of the processing chain which had always be a great advantage of \TeX ; but many didn't mind.

¹³ Before $\text{X}\mathcal{F}\text{\TeX}$ there had been $\text{\TeX}\text{GX}$ (on the Mac only) which used the same series of Apple libraries, then called “GX technology”—TrueType GX has been an extensions of the TrueType font format, now replaced and enhanced by AAT—Apple Advanced Typography.

3.2 Towards the infinite and beyond

Another \TeX engine worth mentioning in passing is called Aleph (\aleph), of Ω -ish ascent. It started as an attempt to stabilize Ω while merging the extensions by $\varepsilon\text{\TeX}$ at the same time (hence its original name, $\varepsilon\text{-}\Omega$).

While it attracted much attention for a few years after it was launched in 2001, it is today shadowed by another successor to \TeX , which is now thought as representing the future path.

4 Howling to the moon

The “successor” is $\text{Lua}\text{\TeX}$. As this seems to be yet another prefixed version of \TeX , we shall first explain what that prefix is. Lua is a small scripting language originated at a university of Rio de Janeiro in Brazil which was developed to be embedded in other applications. The world “lua” means moon in Portuguese (hence the title of that section).

So the idea seems clear, $\text{Lua}\text{\TeX}$ is $\text{Lua} + \text{\TeX}$: an embedded language in \TeX , enabling to go even further than everything that could be done before with macros; in $\text{Lua}\text{\TeX}$ we will also have the Lua language and we can write Lua functions in addition to \TeX macros. But ... there is actually much more: while $\text{Lua}\text{\TeX}$ was indeed meant to be $\text{Lua} + \text{\TeX}$ (actually $\text{pdf}\text{\TeX}$, now merged into $\text{pdf}\varepsilon\text{\TeX}$) when it was first conceived in the beginning of 2005, it is also incorporating the features of \aleph and its Ω predecessor, therefore effectively merging two family of engines: the “ Ω way” and the “ $\text{pdf}\text{\TeX}$ way”. Lua will be present at every stage of the processing chain, with *callbacks* enabling the user to redefine parts of \TeX 's tasks using Lua functions. Finally, METAPOST is planned to be part of it too, being rewritten as a library (instead of a stand-alone program).

$\text{Lua}\text{\TeX}$ is under active development today and is planned to be released at the 2008 TUG in Ireland.¹⁴

5 Back to the future

With $\text{Lua}\text{\TeX}$ we touch the most recent developments in \TeX , and here it seems nice to say some words to summarize the changes that we have seen above.

If any general view is to be had, it seems to me that the main changes that \TeX has undergone over the years were not only major improvements but genuine *Copernician Revolutions* which progressively widened \TeX 's field of application. I have

¹⁴ When TUG will return to Cork which it had already made famous in the \TeX world 18 years ago!

tried to classify those phases in the article by making each one of them a different section: section 1 shows how $\text{T}_{\text{E}}\text{X}$ started with an approach of typesetting akin to that of the craftsman’s carefully setting type to build a page of text,¹⁵ while undergoing an initial modest expansion along with some “standardization” (Babel package for $\text{L}^{\text{A}}\text{T}_{\text{E}}\text{X}$, Cork encoding). Then it went through a phase where the first experimentations were made to handle “complex scripts” (section 2) and this gave birth to the first true extensions of $\text{T}_{\text{E}}\text{X}$ which are actually quite old (again, $\text{T}_{\text{E}}\text{X--X}_{\text{E}}\text{T}$ was written 20 years ago). These extensions were consolidated in the recent past described in section 3, when $\text{T}_{\text{E}}\text{X}$ showed how it could still keep up with the major improvements of the printing industry (PDFformat, TrueType and then Opentype fonts). The present of $\text{T}_{\text{E}}\text{X}$ development, section 4 is exemplified by $\text{LuaT}_{\text{E}}\text{X}$ which, once again, comes with a complete change of perspective on $\text{T}_{\text{E}}\text{X}$ processing. Thus have been the four “ages” of $\text{T}_{\text{E}}\text{X}$.

Writing history is important, and I have always got the impression that the $\text{T}_{\text{E}}\text{X}$ community did not care enough about its own history—there are of course well-known bits and anecdotes about $\text{T}_{\text{E}}\text{X}$, but those are closer to legend than to history. Writing history is also a difficult and humble task, and I cannot claim having covered everything that was worthy remembering about $\text{T}_{\text{E}}\text{X}$: some words could have been said about macro packages (beyond $\text{L}^{\text{A}}\text{T}_{\text{E}}\text{X}$, which I simply quoted in passing) as well as $\text{T}_{\text{E}}\text{X}$ distributions which have contributed a lot to $\text{T}_{\text{E}}\text{X}$ ’s expansion. Nor have I talked about important “industry standards” like XML which have also become an important part of $\text{T}_{\text{E}}\text{X}$ ’s capability today (this would have been more linked to the macro packages than to the engines themselves). It is therefore my hope that we can, maybe, enhance this article with more descriptions and memories, and I have opened a small page at the $\text{ConT}_{\text{E}}\text{Xt}$ wiki to discuss this: http://wiki.contextgarden.net/History_of_TeX.

As a conclusion I would like to name a few places where I’ve encountered $\text{T}_{\text{E}}\text{X}$ personally, as it gives an idea of the versatility of $\text{T}_{\text{E}}\text{X}$ and the extent of the *Pax T_EXnica*:

- The general scientific community and especially mathematicians and computer scientists.¹⁶
- People from humanities, especially in Ancient Greek and linguistics.¹⁷

¹⁵ Let us not forget: “Rhymes are typeset with boxes and glue”, in *The T_EXbook*, chapter 14.

¹⁶ This was obvious but I felt I still had to mention it first!

¹⁷ Indeed, what other free program can handle at the same time Ancient Greek, Russian, Lituianian, Latvian, Sanskrit

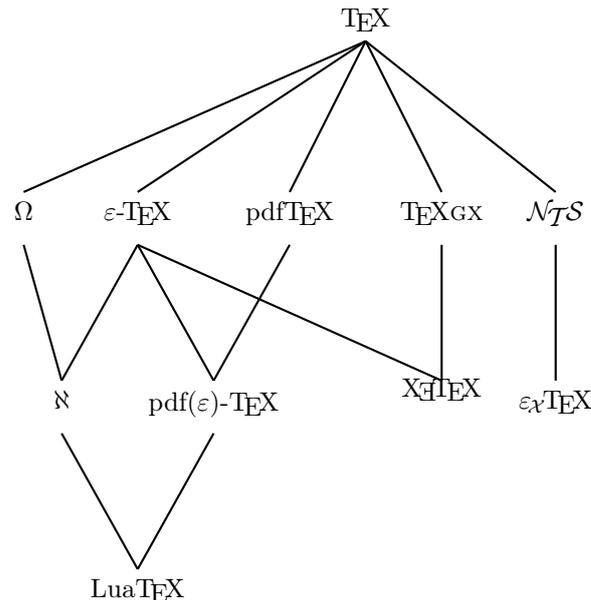


Figure 4: The happy $\text{T}_{\text{E}}\text{X}$ family. The different extensions have been divided into successive “generations” of engine, corresponding to the different sections in this article.

- Shopkeeper from one of the biggest Chinese book-shops in Paris.¹⁸
- Musicians needing to engrave scores.¹⁹
- Users of free software.²⁰
- People involved in the publishing industry.²¹

and French? Someone at my university was doing a Master on Indo-European linguistics and did really need to input all these languages.

¹⁸ Perhaps my most amazing encounter with $\text{T}_{\text{E}}\text{X}$ in a place I didn’t expect it at all, but I swear it is true: while gazing at the shelves of the aforementioned bookshops I overheard two members of the staff discussing how to produce documents in Chinese (probably for the shop’s catalog).

¹⁹ $\text{MusicT}_{\text{E}}\text{X}$ and $\text{MusixT}_{\text{E}}\text{X}$ have many an adept in spite of their extreme difficulty to master.

²⁰ Indeed, on an average Linux distribution, there are very few softwares able to rival with $\text{T}_{\text{E}}\text{X}$ —OpenOffice is an obvious example, but it may be the only other one.

²¹ Especially for processing XML, as already mentioned.