

low level

TEX

characters

Contents

1	Introduction	1
2	History	1
3	The heritage	2
4	The LMTX approach	3

1 Introduction

This explanation is part of the low level manuals because in practice users will not have to deal with these matters in MkIV and even less in LMTX. You can skip to the last section for commands.

2 History

If we travel back in time to when $\text{T}_{\text{E}}\text{X}$ was written we end up in eight bit character universe. In fact, the first versions assumed seven bits, but for comfortable use with languages other than English that was not sufficient. Support for eight bits permits the usage of so called code pages as supported by operating systems. Although ascii input became kind of the standard soon afterwards, the engine can be set up for different encodings. This is not only true for $\text{T}_{\text{E}}\text{X}$, but for many of its companions, like MetaFont and therefore MetaPost.¹

Core components of a $\text{T}_{\text{E}}\text{X}$ engine are hyphenation of words, applying inter-character kerns and build ligatures. In traditional $\text{T}_{\text{E}}\text{X}$ engines those processes are interwoven into the par builder but in Lua $\text{T}_{\text{E}}\text{X}$ these are separate stages. The original approach is the reason that there is a relation between the input encoding and the font encoding: the character in the input is the slot used in a reference to a glyph. When producing the final result (e.g. pdf) there can also be a mapping to an index in a font resource.

input A [tex ->] font slot A [backend ->] glyph index A

The mapping that $\text{T}_{\text{E}}\text{X}$ does is normally one-to-one but an input character can undergo some transformation. For instance a character beyond ascii 126 can be made active and expand to some character number that then becomes the font slot. So, it is the expansion (or meaning) of a character that end up as numeric reference in the glyph

¹ This remapping to an internal representation (e.g. ebcdic) is not present in Lua $\text{T}_{\text{E}}\text{X}$ where we assume utf8 to be the input encoding. The MetaPost library that comes with Lua $\text{T}_{\text{E}}\text{X}$ still has that code but in LuaMeta $\text{T}_{\text{E}}\text{X}$ it's gone. There one can set up the machinery to be utf8 aware too.

node. Virtual fonts can introduce yet another remapping but that's only visible in the backend.

Actually, in Lua \TeX the same happens but in practice there is no need to go active because (at least in Con \TeX t) we assume a Unicode path so there the font slot is the Unicode got from the utf8 input.

In the eight bit universe macro packages (have to) provide all kind of means to deal with (in the perspective of English) special characters. For instance, `\"a` would put a diaeresis on top of the a or even better, refer to a character in the encoding that the chosen font provides. Because there are some limitations of what can go in an eight bit font, and because in different countries the used \TeX fonts evolved kind of independent, we ended up with quite some different variants of fonts. It was only with the Latin Modern project that this became better. Interesting is that when we consider the fact that such a font has often also hardly used symbols (like registered or copyright) coming up with an encoding vector that covers most (latin based) European languages (scripts) is not impossible² Special symbols could simply go into a dedicated font, also because these are always accessed via a macro so who cares about the input. It never happened.

Keep in mind that when utf8 is used with eight bit engines, Con \TeX t will convert sequences of characters into a slot in a font (depending on the font encoding used which itself depends on the coverage needed). For this every first (possible) byte of a multi-byte utf sequence is an active character, which is no big deal because these are outside the ascii range. Normal ascii characters are single byte utf sequences and fall through without treatment.

Anyway, in Con \TeX t MkII we dealt with this by supporting mixed encodings, depending on the (local) language, referencing the relevant font. It permits users to enter the text in their preferred input encoding and also get the words properly hyphenated. But we can leave these MkII details behind.

3 The heritage

In MkIV we got rid of input and font encodings, although one can still load files in a specific code page.³ We also kept the means to enter special characters, if only because text editors seldom support(ed) a wide range of visual editing of those. This is why we still have

² And indeed in the Latin Modern project we came up with one but it was already to late for it to become popular.

³ I'm not sure if users ever depend on an input encoding different from utf8.

```
\"u \^a \v{s} \AE \ij \eacute \oslash
```

and many more. The ones with one character names are rather common in the T_EX community but it is definitely a weird mix of symbols. The next two are kind of outdated: in these days you delegate that to the font handler, where turning them into ‘single’ character references depends on what the font offers, how it is set up with respect to (for instance) ligatures, and even might depend on language or script.

The ones with the long names partly are tradition, but as we have a lot of them, in MkII they actually serve a purpose. These verbose names are used in the so called encoding vectors and are part of the utf expansion vectors. They are also used in labels so that we have a good indication if what goes in there: remember that in those times editors often didn't show characters, unless the font for display had them, or the operating system somehow provided them from another font. These verbose names are used for latin, greek and cyrillic and for some other scripts and symbols. They take up quite a bit of hash space and the format file.⁴

4 The LMTX approach

In the process of tagging all (public) macros in LMTX (which happened in 2020-2021) I wondered if we should keep these one character macros, the references to special characters and the verbose ones. When asked on the mailing list it became clear that users still expect the short ones to be present, often just because old bibT_EX files are used that might need them. However, in MkIV and LMTX we load bibT_EX files in a way that turn these special character references into proper utf8 input so it makes a weak argument. Anyway, although they could go, for now we keep them because users expect them. However, in LMTX the implementation is somewhat different now, a bit more efficient in terms of hash and memory, potentially a bit less efficient in runtime, but no one will notice that.

A new command has been introduced, the very short `\chr`.

```
\chr {a`} \chr {a´} \chr {a¨}
\chr {`a} \chr {'a} \chr {"a}
\chr {a acute} \chr {a grave} \chr {a umlaut}
\chr {aacute} \chr {agrave} \chr {aumlaut}
```

In the first line the composed character using two characters, a base and a so called mark. Actually, one doesn't have to use `\chr` in that case because ConT_EXt does already

⁴ In MkII we have an abstract front-end with respect to encodings and also an abstract backend with respect to supported drivers but both approaches no longer make sense today.

collapse characters for you. The second line looks like the shortcuts `\``, `\'` and `\"`. The third and fourth lines could eventually replace the more symbolic long names, if we feel the need. Watch out: in Unicode input the marks come *after*.

```
à á ä
à á ä
á à a~mla~t
á à a~mla~t
```

Currently the repertoire is somewhat limited but it can be easily be extended. It all depends on user needs (doing Greek and Cyrillic for instance). The reason why we actually save code deep down is that the helpers for this have always been there.⁵

The `\` commands are now just aliases to more verbose and less hackery looking macros:

```
\withgrave      à  \`  à
\withacute      á  \'  á
\withcircumflex â  \^  â
\withtilde      ã  \~  ã
\withmacron     ā  \=  ā
\withbreve      ě  \u  ě
\withdotaccent  ċ  \.  .c
\withdiaeresis  ë  \"  ë
\withring       ũ  \r  ũ
\withhungarumlaut  ű  \H  ű
\withcaron      ě  \v  ě
\withcedilla    ç  \c  ç
\withogonek     ę  \k  ę
```

Not all fonts have these special characters. Most natural is to have them available as precomposed single glyphs, but it can be that they are just two shapes with the marks anchored to the base. It can even be that the font somehow overlays them, assuming (roughly) equal widths. The compose font feature in ConT_EXt normally can handle most well.

An occasional ugly rendering doesn't matter that much: better have something than nothing. But when it's the main language (script) that needs them you'd better look for a font that handles them. When in doubt, in ConT_EXt you can enable checking:

⁵ So if needed I can port this approach back to MkIV, but for now we keep it as is because we then have a reference.

command	equivalent to
<code>\checkmissingcharacters</code>	<code>\enabletrackers[fonts.missing]</code>
<code>\removemissingcharacters</code>	<code>\enabletrackers[fonts.missing=remove]</code>
<code>\replacemissingcharacters</code>	<code>\enabletrackers[fonts.missing=replace]</code>
<code>\handlemissingcharacters</code>	<code>\enabletrackers[fonts.missing={decompose,replace}]</code>

The decompose variant will try to turn a composed character into its components so that at least you get something. If that fails it will inject a replacement symbol that stands out so that you can check it. The console also mentions missing glyphs. You don't need to enable this by default⁶ but you might occasionally do it when you use a font for the first time.

In LMTX this mechanism has been upgraded so that replacements follow the shape and are actually real characters. The decomposition has not yet been ported back to MkIV.

The full list of commands can be queried when a tracing module is loaded:

```
\usemodule[s][characters-combinations]
```

```
\showcharactercombinations
```

We get this list:

acute	U+00301	´	<code>\withacute</code>
breve	U+00306	˘	<code>\withbreve</code>
caron	U+0030C	ˇ	<code>\withcaron</code>
caron below	U+0032C	̣	<code>\withcaronbelow</code>
cedilla	U+00327	¸	<code>\withcedilla</code>
circumflex	U+00302	ˆ	<code>\withcircumflex</code>
circumflex below	U+0032D	̂	<code>\withcircumflexbelow</code>
comma below	U+00327	¸	<code>\withcommabelow</code>
diaeresis	U+00308	¨	<code>\withdiaeresis</code>
dieresis	U+00308	¨	<code>\withdieresis</code>
dot	U+00307	·	<code>\withdot</code>
dot below	U+00323	̣	<code>\withdotbelow</code>
double acute	U+0030B	˚	<code>\withdoubleacute</code>
double grave	U+0030F	˘	<code>\withdoublegrave</code>
double vertical line	U+0030E	="	<code>\withdoubleverticalline</code>
grave	U+00300	˘	<code>\withgrave</code>
hook	U+00309	ˆ	<code>\withhook</code>

⁶ There is some overhead involved here.

hook below	U+1FA9D		\withhookbelow
hungarumlaut	U+0030B	“	\withhungarumlaut
inverted breve	U+00311	^	\withinvertedbreve
line	U+00304	-	\withline
line below	U+00331	_	\withlinebelow
macron	U+00304	-	\withmacron
macron below	U+00331	_	\withmacronbelow
middle dot	U+000B7	·	\withmiddledot
ogonek	U+00328	˘	\withogonek
overline	U+00305	—	
ring	U+0030A	°	\withring
ring below	U+00325	◦	\withringbelow
slash	U+0002F	/	\withslash
stroke	U+0002F	/	\withstroke
tilde	U+00303	~	\withtilde
tilde below	U+00330	˜	\withtildebelow
vertical line	U+0030D	¦	\withverticalline

Some combinations are special for ConT_EXt because Unicode doesn't specify decomposition for all composed characters.

4 Colofon

Author	Hans Hagen
ConT _E Xt	2023.04.27 17:04
LuaMetaT _E X	2.1008
Support	www.pragma-ade.com contextgarden.net