

Providing some UTF-8 support via `inputenc`

Frank Mittelbach Chris Rowley*

? 2015/06/27 v1.1n UTF-8 support for inputenc? printed July 7,
2015

This file is maintained by the L^AT_EX Project team.
Bug reports can be opened (category `latex`) at
<http://latex-project.org/bugs.html>.

Contents

1	Introduction	2
1.1	Background and general stuff	2
1.2	More specific stuff	2
1.3	Notes	3
1.4	Basic operation of the code	3
2	Coding	4
2.1	Housekeeping	4
2.2	Parsing UTF-8 input	4
2.3	Mapping Unicode codes to L ^A T _E X internal forms	6
2.4	Loading Unicode mappings at begin document	8
3	Mapping characters — based on font (glyph) encodings	10
3.1	About the table itself	10
3.2	The mapping table	10
3.3	Notes	20
3.4	Mappings for OT1 glyphs	20
3.5	Mappings for OMS glyphs	21
3.6	Mappings for TS1 glyphs	21
3.7	Mappings for <code>latex.ltx</code> glyphs	21
4	A test document	22

*Borrowing heavily from code by David Carlisle and tables by Sebastian Rahtz; some table and code cleanup by Javier Bezos

1 Introduction

[The whole section is rather unfinished ... just like the code, sorry!]

1.1 Background and general stuff

For many reasons what this package provides is a long way from any type of ‘Unicode compliance’.

In stark contrast to 8-bit character sets, with 16 or more bits it can easily be very inefficient to support the full range.¹ Moreover, useful support of character input by a typesetting system overwhelmingly means finding an acceptable visual representation of a sequence of characters and this, for L^AT_EX, means having available a suitably encoded 8-bit font.

Unfortunately it is not possible to predict exactly what valid UTF-8 octet sequences will appear in a particular file so it is best to make all the unsupported but valid sequences produce a reasonably clear and noticeable error message.

There are two directions from which to approach the question of what to load. One is to specify the ranges of Unicode characters that will result in some sensible typesetting; this requires the provider to ensure that suitable fonts are loaded and that these input characters generate the correct typesetting via the encodings of those fonts. The other is to inspect the font encodings to be used and use these to define which input Unicode characters should be supported.

For Western European languages, at least, going in either direction leads to many straightforward decisions and a few that are more subjective. In both cases some of the specifications are T_EX specific whilst most are independent of the particular typesetting software in use.

As we have argued elsewhere, L^AT_EX needs to refer to characters via ‘seven-bit-text’ names and, so far, these have been chosen by reference to historical sources such as Plain T_EX or Adobe encoding descriptions. It is unclear whether this ad hoc naming structure should simply be extended or whether it would be useful to supplement it with standardised internal Unicode character names such as one or more of the following:²

```
\ltxutwochar <4 hex digits>

\ltxuchar {<hex digits>}
  B H U R R R

\ltxueightchartwo <2 utf8 octets as 8-bit char tokens>
\ltxueightcharthree <3 utf8 octets ...>
\ltxueightcharfour <4 utf8 octets ...>
```

1.2 More specific stuff

In addition to setting up the mechanism for reading UTF-8 characters and specifying the L^AT_EX-level support available, this package contains support for some

¹In fact, L^AT_EX’s current 8-bit support does not go so far as to make all 8-bit characters into valid input.

²Burkhard und Holger Mittelbach spielen mit mir! Sie haben etwas hier geschrieben.

default historically expected \TeX -related characters and some example ‘Unicode definition files’ for standard font encodings.

1.3 Notes

This package does not support Unicode combining characters as \TeX is not really equipped to make this possible.

No attempt is made to be useful beyond Latin, and maybe Cyrillic, for European languages (as of now).

1.4 Basic operation of the code

The `inputenc` package makes the upper 8-bit characters active and assigns to all of them an error message. It then waits for the input encoding files to change this set-up. Similarly, whenever `\inputencoding` is encountered in a document, first the upper 8-bit characters are set back to produce an error and then the definitions for the new input encoding are loaded, changing some of the previous settings.

The 8-bit input encodings currently supported by `inputenc` all use declarations such as `\DeclareInputText` and the like to map an 8-bit number to some \LaTeX internal form, e.g. to `\"a`.

The situation when supporting UTF-8 as the input encoding is different, however. Here we only have to set up the actions of those 8-bit numbers that can be the first octet in a UTF-8 representation of a Unicode character. But we cannot simply set this to some internal \LaTeX form since the Unicode character consists of more than one octet; instead we have to define this starting octet to parse the right number of further octets that together form the UTF-8 representation of some Unicode character.

Therefore when switching to `utf8` within the `inputenc` framework the characters with numbers (hex) from `"C2` to `"DF` are defined to parse for a second octet following, the characters from `"E0` to `"EF` are defined to parse for two more octets and finally the characters from `"F0` to `"F3` are defined to parse for three additional octets. These additional octets are always in the range `"80` to `"B9`.

Thus, when such a character is encountered in the document (so long as expansion is not prohibited) a defined number of additional octets (8-bit characters) are read and from them a unique control sequence name is immediately constructed.

This control sequence is either defined (good) or undefined (likely); in the latter case the user gets an error message saying that this UTF-8 sequence (or, better, Unicode character) is not supported.

If the control sequence is set up to do something useful then it will expand to a \LaTeX internal form: e.g. for the `utf8` sequence of two octets `"C3 "A4` we get `\"a` as the internal form which then, depending on the font encoding, eventually resolves to the single glyph ‘latin-a-umlaut’ or to the composite glyph ‘latin-a with an umlaut accent’.

These mappings from (UTF-8 encoded) Unicode characters to \LaTeX internal forms are made indirectly. The code below provides a declaration `\DeclareUnicodeCharacter` which maps Unicode numbers (as hexadecimal) to \LaTeX internal forms.

This mapping needs to be set up only once so it is done at `\begin{document}` by looking at the list of font encodings that are loaded by the document and providing mappings related to those font encodings whenever these are available.

Thus at most only those Unicode characters that can be represented by the glyphs available in these encodings will be defined.

Technically this is done by loading one file per encoding, if available, that is supposed to provide the necessary mapping information.

2 Coding

2.1 Housekeeping

The usual introductory bits and pieces:

```

1 <utf8>\ProvidesFile{utf8.def}
2 <test>\ProvidesFile{utf8-test.tex}
3 <+lcy> \ProvidesFile{lcyenc.dfu}
4 <+ly1> \ProvidesFile{ly1enc.dfu}
5 <+oms> \ProvidesFile{omsenc.dfu}
6 <+ot1> \ProvidesFile{ot1enc.dfu}
7 <+ot2> \ProvidesFile{ot2enc.dfu}
8 <+t1> \ProvidesFile{t1enc.dfu}
9 <+t2a> \ProvidesFile{t2aenc.dfu}
10 <+t2b> \ProvidesFile{t2benc.dfu}
11 <+t2c> \ProvidesFile{t2cenc.dfu}
12 <+ts1> \ProvidesFile{ts1enc.dfu}
13 <+x2> \ProvidesFile{x2enc.dfu}
14 <+all> \ProvidesFile{utf8enc.dfu}
15 [2015/06/27 v1.1n UTF-8 support for inputenc]

16 <*utf8>
17 \makeatletter

```

We restore the `\catcode` of space (which is set to ignore in `inputenc`) while reading `.def` files. Otherwise we would need to explicitly use `\space` all over the place in error and log messages.

```

18 \catcode'\ \saved@space@catcode

```

2.2 Parsing UTF-8 input

`\UTFviii@two@octets`
`\UTFviii@three@octets`
`\UTFviii@four@octets`

A UTF-8 char (that is not actually a 7-bit char, i.e. a single octet) is parsed as follows: each starting octet is an active `TEX` character token; each of these is defined below to be a macro with one to three arguments nominally (depending on the starting octet). It calls one of `\UTFviii@two@octets`, `\UTFviii@three@octets`, or `\UTFviii@four@octets` which then actually picks up the argument(s).

From the arguments a control sequence with a name of the form `u8:#1#2...` is constructed where the `#i` ($i > 1$) are the arguments and `#1` is the starting octet (as a `TEX` character token). Since some or even all of these characters are active (when `inputenc` is loaded) we need to use `\string` when building the `csname`.

The `csname` thus constructed can of course be undefined but to avoid producing an unhelpful low-level undefined command error we pass it to `\UTFviii@defined` which is responsible for producing a more sensible error message (not yet done!!). If, however, it is defined we simply execute the thing (which should then expand to an encoding specific internal L^AT_EX form).

```

19 \def\UTFviii@two@octets#1#2{\expandafter
20   \UTFviii@defined\csname u8:#1\string#2\endcsname}

```

```

21 \def\UTFviii@three@octets#1#2#3{\expandafter
22   \UTFviii@defined\csname u8:#1\string#2\string#3\endcsname}
23 \def\UTFviii@four@octets#1#2#3#4{\expandafter
24   \UTFviii@defined\csname u8:#1\string#2\string#3\string#4\endcsname}
\UTFviii@defined This tests whether its argument is different from \relax: it either calls for a
sensible error message (not done), or it gets the \fi out of the way (in case the
command has arguments) and executes it.
25 \def\UTFviii@defined#1{%
26   \ifx#1\relax
The endlene character has a special definition within the inputenc package (it is
gobbling spaces). For this reason we can't produce multiline strings without some
precaution.
27     \PackageError{inputenc}{Unicode\space char\space \string#1\space
28                               not\space set\space up\space
29                               for\space use\space with\space LaTeX}\@eha
30   \else\expandafter
31     #1%
32   \fi
33 }
\UTFviii@loop This wonderful bit of code from Dr Carlisle defines the starting octets to call
\UTFviii@two@octets etc as appropriate. The starting octet itself is passed di-
rectly as the first argument, the others are picked up later en route.
The \UTFviii@loop loops through the numbers starting at \count@ and end-
ing at \@tempcnta - 1, each time executing the code in \UTFviii@tmp.
All this is done in a group so that temporary catcode changes etc. vanish after
everything is set up.
It may be a good idea to add code to deal with 'illegal utf8 octets': at present
these will be handled by whatever code was in use for 8-bit input before this code
is executed.
34 \begingroup
35 \catcode'\~13
36 \catcode'\~12
37 \def\UTFviii@loop{%
38   \uccode'\~\count@
39   \uppercase\expandafter{\UTFviii@tmp}%
40   \advance\count@\@ne
41   \ifnum\count@<\@tempcnta
42     \expandafter\UTFviii@loop
43   \fi}
Setting up 2-byte UTF-8:
44   \count@"C2
45   \@tempcnta"E0
46   \def\UTFviii@tmp{\xdef~{\noexpand\UTFviii@two@octets\string~}}
47 \UTFviii@loop
Setting up 3-byte UTF-8:
48   \count@"E0
49   \@tempcnta"F0
50   \def\UTFviii@tmp{\xdef~{\noexpand\UTFviii@three@octets\string~}}
51 \UTFviii@loop

```

Setting up 4-byte UTF-8:

```
52 \count@"F0
53 \@tempcnta"F4
54 \def\UTFviii@tmp{\xdef~{\noexpand\UTFviii@four@octets\string~}}
55 \UTFviii@loop
56 \endgroup
```

For this case we must disable the warning generated by `inputenc` if it doesn't see any new `\DeclareInputText` commands.

```
57 \@inpenc@test
```

If this file (`utf8.def`) is not being read while setting up `inputenc`, i.e. in the preamble, but when `\inputencoding` is called somewhere within the document, we do not need to input the specific Unicode mappings again. We therefore stop reading the file at this point.

```
58 \ifx\@begindocumenthook\@undefined
59 \makeatother
```

The `\fi` must be on the same line as `\endinput` or else it will never be seen!

```
60 \endinput \fi
```

2.3 Mapping Unicode codes to L^AT_EX internal forms

`\DeclareUnicodeCharacter` The `\DeclareUnicodeCharacter` declaration defines a mapping from a Unicode character code point to a L^AT_EX internal form. The first argument is the Unicode number as hexadecimal digits and the second is the actual L^AT_EX internal form.

We start by making sure that some characters have the right `\catcode` when they are used in the definitions below.

```
61 \begingroup
62 \catcode'\="=12
63 \catcode'\<=12
64 \catcode'\.=12
65 \catcode'\,=12
66 \catcode'\;=12
67 \catcode'\!=12
68 \catcode'\~=13

69 \gdef\DeclareUnicodeCharacter#1#2{%
70 \count@"#1\relax
71 \wlog{ \space\space defining Unicode char U+#1 (decimal \the\count@)}%
72 \begingroup
```

Next we do the parsing of the number stored in `\count@` and assign the result to `\UTFviii@tmp`. Actually all this could be done in-line, the macro `\parse@XML@charref` is only there to extend this code to parsing Unicode numbers in other contexts one day (perhaps).

```
73 \parse@XML@charref
```

Here is an example of what is happening, for the pair "C2 "A3 (which is the utf8 representation for the character £). After `\parse@XML@charref` we have, stored in `\UTFviii@tmp`, a single command with two character tokens as arguments:

```
[tC2 and tA3 are the characters corresponding to these two octets]
\UTFviii@two@octets tC2tA3
```

what we actually need to produce is a definition of the form

```
\def\u8:tC2tA3 {\LATEX internal form}.
```

So here we temporarily redefine the prefix commands `\UTFviii@two@octets`, etc. to generate the csname that we wish to define; the `\strings` are added in case these tokens are still active.

```
74 \def\UTFviii@two@octets##1##2{\csname u8:##1\string##2\endcsname}%
75 \def\UTFviii@three@octets##1##2##3{\csname u8:##1%
76 \string##2\string##3\endcsname}%
77 \def\UTFviii@four@octets##1##2##3##4{\csname u8:##1%
78 \string##2\string##3\string##4\endcsname}%
```

Now we simply:-) need to use the right number of `\expandafters` to finally construct the definition: expanding `\UTFviii@tmp` once to get its contents, a second time to replace the prefix command by its `\csname` expansion, and a third time to turn the expansion into a csname after which the `\gdef` finally gets applied. We add an irrelevant `\IeC` and braces around the definition, in order to avoid any space after the command being gobbled up when the text is written out to an auxiliary file (see `inputenc` for further details

```
79 \expandafter\expandafter\expandafter
80 \expandafter\expandafter\expandafter
81 \expandafter
82 \gdef\UTFviii@tmp{\IeC{#2}}%
83 \endgroup
84 }
```

`\parse@XML@charref` This macro parses a Unicode number (decimal) and returns its UTF-8 representation as a sequence of non-active T_EX character tokens. In the original code it had two arguments delimited by `;` here, however, we supply the Unicode number implicitly.

```
85 \gdef\parse@XML@charref{%
```

We need to keep a few things local, mainly the `\uccode`'s that are set up below. However, the group originally used here is actually unnecessary since we call this macro only within another group; but it will be important to restore the group if this macro gets used for other purposes.

```
86 % \begingroup
```

The original code from David supported the convention that a Unicode slot number could be given either as a decimal or as a hexadecimal (by starting with `x`). We do not do this so this code is also removed. This could be reactivated if one wants to support document commands that accept Unicode numbers (but then the first case needs to be changed from an error message back to something more useful again).

```
87 % \uppercase{\count@{if x\noexpand#1"\else#1\fi#2}}\relax
```

As `\count@` already contains the right value we make `\parse@XML@charref` work without arguments.

```
88 \ifnum\count@<"A0\relax
89 \PackageError{inputenc}{Cannot\space define\space Unicode\space
90 char\space value\space <\space 00A0}\@eha
```

Do not ask us to provide an explanation for the code below, it is borrowed straight from `xmltex` by David and we trust him totally (and we are too lazy to reread the Unicode book to see if this is the correct algorithm).³

```

91 \else\ifnum\count@<"800\relax
92   \parse@UTFviii@a,%
93   \parse@UTFviii@b C\UTFviii@two@octets.%,%
94 \else\ifnum\count@<"10000\relax
95   \parse@UTFviii@a;%
96   \parse@UTFviii@a,%
97   \parse@UTFviii@b E\UTFviii@three@octets.{,;}%
98 \else
99   \parse@UTFviii@a;%
100  \parse@UTFviii@a,%
101  \parse@UTFviii@a!%
102  \parse@UTFviii@b F\UTFviii@four@octets.{!,;}%
103  \fi
104  \fi
105  \fi
106 % \endgroup
107 }

```

`\parse@UTFviii@a` ...so somebody else can document this part :-) ... David?:-))))!

```

108 \gdef\parse@UTFviii@a#1{%
109   \@tempcnta\count@
110   \divide\count@ 64
111   \@tempcntb\count@
112   \multiply\count@ 64
113   \advance\@tempcnta-\count@
114   \advance\@tempcnta 128
115   \uccode'#1\@tempcnta
116   \count@\@tempcntb}

```

`\parse@UTFviii@b` ... same here

```

117 \gdef\parse@UTFviii@b#1#2#3#4{%
118   \advance\count@ "#10\relax
119   \uccode'#3\count@
120   \uppercase{\gdef\UTFviii@tmp{#2#3#4}}
121 \endgroup

```

```
122 \@onlypreamble\DeclareUnicodeCharacter
```

These are preamble only as long as we don't support Unicode charrefs in documents.

```

123 \@onlypreamble\parse@XML@charref
124 \@onlypreamble\parse@UTFviii@a
125 \@onlypreamble\parse@UTFviii@b

```

2.4 Loading Unicode mappings at begin document

The original plan was to set up the UTF-8 support at `\begin{document}`; but then any text characters used in the preamble (as people do even though advised

³We were hoping to also find in his work the T_EX code for going the other way: from UTF-8 octets to Unicode slot number, but no luck!

against it) would fail in one way or the other. So the implementation was changed and the Unicode definition files for already defined encodings are loaded here.

We loop through all defined font encodings (stored in `\cdp@list`) and for each load a file `nameenc.dfu` if it exist. That file is then supposed to contain `\DeclareUnicodeCharacter` declarations.

```

126 \begingroup
127 \def\cdp@elt#1#2#3#4{%
128   \wlog{Now handling font encoding #1 ...}%
129   \lowercase{%
130     \InputIfFileExists{#1enc.dfu}}%
131     {\wlog{... processing UTF-8 mapping file for font %
132       encoding #1}}%

```

The previous line is written to the log with the newline char being ignored (thus not producing a space). Therefore either everything has to be on a single input line or some special care must be taken. From this point on we ignore spaces again, i.e., while we are reading the `.dfu` file. The `\endgroup` below will restore it again.

```

133       \catcode'\ 9\relax}%
134       {\wlog{... no UTF-8 mapping file for font encoding #1}}%
135   }
136 \cdp@list
137 \endgroup

```

However, we don't know if there are font encodings still to be loaded (either with `fontenc` or directly with `\input` by some some package). Font encoding files are loaded only if the corresponding encoding has not been loaded yet, and they always begin with `\DeclareFontEncoding`. We now redefine the internal kernel version of the latter to load the Unicode file if available.

```

138 \def\DeclareFontEncoding@#1#2#3{%
139   \expandafter
140   \ifx\csname T@#1\endcsname\relax
141     \def\cdp@elt{\noexpand\cdp@elt}%
142     \xdef\cdp@list{\cdp@list\cdp@elt{#1}%
143       {\default@family}{\default@series}%
144       {\default@shape}}%
145     \expandafter\let\csname#1-cmd\endcsname\@changed@cmd
146     \begingroup
147       \wlog{Now handling font encoding #1 ...}%
148       \lowercase{%
149         \InputIfFileExists{#1enc.dfu}}%
150         {\wlog{... processing UTF-8 mapping file for font %
151           encoding #1}}%
152         {\wlog{... no UTF-8 mapping file for font encoding #1}}%
153       \endgroup
154   \else
155     \@font@info{Redeclaring font encoding #1}%
156   \fi
157   \global\@namedef{T@#1}{#2}%
158   \global\@namedef{M@#1}{\default@M#3}%
159   \xdef\LastDeclaredEncoding{#1}%
160   }
161 \endgroup

```

3 Mapping characters — based on font (glyph) encodings

This section is a first attempt to provide Unicode definitions for characters whose standard glyphs are currently provided by the standard \LaTeX font-encodings `T1`, `OT1`, etc. They are by no means completed and need checking.

For example, one should check the already existing input encodings for glyphs that may in fact be available and required, e.g. `latin4` has a number of glyphs with the `\=` accent. Since the `T1` encoding does not provide such glyphs, these characters are not listed below (yet).

The list below was generated by looking at the current \LaTeX font encoding files, e.g., `t1enc.def` and using the work by Sebastian Rahtz (in `ucharacters.sty`) with a few modifications. In combinations such as `\^i` the preferred form is that and not `\~i`.

This list has been built from several sources, obviously including the Unicode Standard itself. These sources include Passive \TeX by Sebastian Rahtz, the `unicode` package by Dominique P. G. Unruh (mainly for Latin encodings) and `text4ht` by Eitan Gurari (for Cyrillic ones).

Note that it strictly follows the Mittelbach principles for input character encodings: thus it offers no support for using utf8 representations of math symbols such as \times or \div (in math mode).

3.1 About the table itself

In addition to generating individual files, the table below is, at present, a one-one (we think) partial relationship between the (ill-defined) set of LICRs and the Unicode slots "0080 to "FFFF. At present these entries are used only to define a collection of partial mappings from Unicode slots to LICRs; each of these mappings becomes full if we add an exception value ('not defined') to the set of LICRs.

It is probably not essential for the relationship in the full table to be one-one; this raises questions such as: the exact role of LICRs; the formal relationships on the set of LICRs; the (non-mathematical) relationship between LICRs and Unicode (which has its own somewhat fuzzy equivalences); and ultimately what a character is and what a character representation and/or name is.

Viewed this a way, the result has, perhaps puzzling, just two (we think) gaps in the second 128 'Unicode slots' (00A0 and 00AD): neither of these is really a character, of course.

It is unclear the extent to which entries in this table should resemble the closely related ones in the 8-bit `inputenc` files. The Unicode standard claims that the first 256 slots 'are' ASCII and Latin-1.

Of course, \TeX itself typically does not treat even many perfectly 'normal text' 7-bit slots as text characters, so it is unclear whether \LaTeX should even attempt to deal in any consistent way with those Unicode slots that are not definitive text characters.

3.2 The mapping table

Note that the first argument must be a 4-hex-digit number greater than 00BF.

There are few notes about inconsistencies etc at the end of the table.

```

162 <all, t1, ot1, ly1>\DeclareUnicodeCharacter{00A1}{\textexclamdown}
163 <all, ts1, ly1>\DeclareUnicodeCharacter{00A2}{\textcent}
164 <all, ts1, t1, ot1, ly1>\DeclareUnicodeCharacter{00A3}{\textsterling}
165 <all, x2, ts1, t2c, t2b, t2a, ly1, lcy>\DeclareUnicodeCharacter{00A4}{\textcurrency}
166 <all, ts1, ly1>\DeclareUnicodeCharacter{00A5}{\textyen}
167 <all, ts1, ly1>\DeclareUnicodeCharacter{00A6}{\textbrokenbar}
168 <all, x2, ts1, t2c, t2b, t2a, oms, ly1>\DeclareUnicodeCharacter{00A7}{\textsection}
169 <all, ts1>\DeclareUnicodeCharacter{00A8}{\textasciidieresis}
170 <all, ts1, utf8>\DeclareUnicodeCharacter{00A9}{\textcopyright}
171 <all, ts1, ly1, utf8>\DeclareUnicodeCharacter{00AA}{\textordfeminine}
172 *all, x2, t2c, t2b, t2a, t1, ot2, ly1, lcy>
173 \DeclareUnicodeCharacter{00AB}{\guillemotleft}
174 /all, x2, t2c, t2b, t2a, t1, ot2, ly1, lcy>
175 <all, ts1>\DeclareUnicodeCharacter{00AC}{\textlnot}
176 <all, ts1, ly1, utf8>\DeclareUnicodeCharacter{00AE}{\textregistered}
177 <all, ts1>\DeclareUnicodeCharacter{00AF}{\textasciimacron}
178 <all, ts1, ly1>\DeclareUnicodeCharacter{00B0}{\textdegree}
179 <all, ts1>\DeclareUnicodeCharacter{00B1}{\textpm}
180 <all, ts1>\DeclareUnicodeCharacter{00B2}{\texttwosuperior}
181 <all, ts1>\DeclareUnicodeCharacter{00B3}{\textthreesuperior}
182 <all, ts1>\DeclareUnicodeCharacter{00B4}{\textasciiacute}
183 <all, ts1, ly1>\DeclareUnicodeCharacter{00B5}{\textmu} % micro sign
184 <all, ts1, oms, ly1>\DeclareUnicodeCharacter{00B6}{\textparagraph}
185 <all, oms, ts1, ly1>\DeclareUnicodeCharacter{00B7}{\textperiodcentered}
186 <all, ot1>\DeclareUnicodeCharacter{00B8}{\c\ }
187 <all, ts1>\DeclareUnicodeCharacter{00B9}{\textonesuperior}
188 <all, ts1, ly1, utf8>\DeclareUnicodeCharacter{00BA}{\textordmasculine}
189 *all, x2, t2c, t2b, t2a, t1, ot2, ly1, lcy>
190 \DeclareUnicodeCharacter{00BB}{\guillemotright}
191 /all, x2, t2c, t2b, t2a, t1, ot2, ly1, lcy>
192 <all, ts1, ly1>\DeclareUnicodeCharacter{00BC}{\textonequarter}
193 <all, ts1, ly1>\DeclareUnicodeCharacter{00BD}{\textonehalf}
194 <all, ts1, ly1>\DeclareUnicodeCharacter{00BE}{\textthreequarters}
195 <all, t1, ot1, ly1>\DeclareUnicodeCharacter{00BF}{\textquestiondown}
196 <all, t1, ly1>\DeclareUnicodeCharacter{00C0}{\@tabacckludge'A}
197 <all, t1, ly1>\DeclareUnicodeCharacter{00C1}{\@tabacckludge'A}
198 <all, t1, ly1>\DeclareUnicodeCharacter{00C2}{\^A}
199 <all, t1, ly1>\DeclareUnicodeCharacter{00C3}{\~A}
200 <all, t1, ly1>\DeclareUnicodeCharacter{00C4}{\"A}
201 <all, t1, ot1, ly1>\DeclareUnicodeCharacter{00C5}{\r A}
202 <all, t1, ot1, ly1, lcy>\DeclareUnicodeCharacter{00C6}{\AE}
203 <all, t1, ly1>\DeclareUnicodeCharacter{00C7}{\c C}
204 <all, t1, ly1>\DeclareUnicodeCharacter{00C8}{\@tabacckludge'E}
205 <all, t1, ly1>\DeclareUnicodeCharacter{00C9}{\@tabacckludge'E}
206 <all, t1, ly1>\DeclareUnicodeCharacter{00CA}{\^E}
207 <all, t1, ly1>\DeclareUnicodeCharacter{00CB}{\"E}
208 <all, t1, ly1>\DeclareUnicodeCharacter{00CC}{\@tabacckludge'I}
209 <all, t1, ly1>\DeclareUnicodeCharacter{00CD}{\@tabacckludge'I}
210 <all, t1, ly1>\DeclareUnicodeCharacter{00CE}{\^I}
211 <all, t1, ly1>\DeclareUnicodeCharacter{00CF}{\"I}
212 <all, t1, ly1>\DeclareUnicodeCharacter{00D0}{\DH}
213 <all, t1, ly1>\DeclareUnicodeCharacter{00D1}{\~N}
214 <all, t1, ly1>\DeclareUnicodeCharacter{00D2}{\@tabacckludge'O}
215 <all, t1, ly1>\DeclareUnicodeCharacter{00D3}{\@tabacckludge'O}

```

```

216 <all,t1,ly1>\DeclareUnicodeCharacter{00D4}{\^O}
217 <all,t1,ly1>\DeclareUnicodeCharacter{00D5}{\~O}
218 <all,t1,ly1>\DeclareUnicodeCharacter{00D6}{\"O}
219 <all,ts1>\DeclareUnicodeCharacter{00D7}{\texttimes}
220 <all,t1,ot1,ly1,lcY>\DeclareUnicodeCharacter{00D8}{\O}
221 <all,t1,ly1>\DeclareUnicodeCharacter{00D9}{\@tabacckludge'U}
222 <all,t1,ly1>\DeclareUnicodeCharacter{00DA}{\@tabacckludge'U}
223 <all,t1,ly1>\DeclareUnicodeCharacter{00DB}{\^U}
224 <all,t1,ly1>\DeclareUnicodeCharacter{00DC}{\"U}
225 <all,t1,ly1>\DeclareUnicodeCharacter{00DD}{\@tabacckludge'Y}
226 <all,t1,ly1>\DeclareUnicodeCharacter{00DE}{\TH}
227 <all,t1,ot1,ly1,lcY>\DeclareUnicodeCharacter{00DF}{\ss}
228 <all,t1,ly1>\DeclareUnicodeCharacter{00E0}{\@tabacckludge'a}
229 <all,t1,ly1>\DeclareUnicodeCharacter{00E1}{\@tabacckludge'a}
230 <all,t1,ly1>\DeclareUnicodeCharacter{00E2}{\^a}
231 <all,t1,ly1>\DeclareUnicodeCharacter{00E3}{\~a}
232 <all,t1,ly1>\DeclareUnicodeCharacter{00E4}{\"a}
233 <all,t1,ly1>\DeclareUnicodeCharacter{00E5}{\r a}
234 <all,t1,ot1,ly1,lcY>\DeclareUnicodeCharacter{00E6}{\ae}
235 <all,t1,ly1>\DeclareUnicodeCharacter{00E7}{\c c}
236 <all,t1,ly1>\DeclareUnicodeCharacter{00E8}{\@tabacckludge'e}
237 <all,t1,ly1>\DeclareUnicodeCharacter{00E9}{\@tabacckludge'e}
238 <all,t1,ly1>\DeclareUnicodeCharacter{00EA}{\^e}
239 <all,t1,ly1>\DeclareUnicodeCharacter{00EB}{\~e}
240 <all,t1,ot1,ly1>\DeclareUnicodeCharacter{00EC}{\@tabacckludge'i}
241 <all,t1,ot1,ly1>\DeclareUnicodeCharacter{00ED}{\@tabacckludge'i}
242 <all,t1,ot1,ly1>\DeclareUnicodeCharacter{00EE}{\^i}
243 <all,t1,ot1,ly1>\DeclareUnicodeCharacter{00EF}{\~i}
244 <all,t1,ly1>\DeclareUnicodeCharacter{00F0}{\dh}
245 <all,t1,ly1>\DeclareUnicodeCharacter{00F1}{\~n}
246 <all,t1,ly1>\DeclareUnicodeCharacter{00F2}{\@tabacckludge'o}
247 <all,t1,ly1>\DeclareUnicodeCharacter{00F3}{\@tabacckludge'o}
248 <all,t1,ly1>\DeclareUnicodeCharacter{00F4}{\^o}
249 <all,t1,ly1>\DeclareUnicodeCharacter{00F5}{\~o}
250 <all,t1,ly1>\DeclareUnicodeCharacter{00F6}{\"o}
251 <all,ts1>\DeclareUnicodeCharacter{00F7}{\textdiv}
252 <all,t1,ot1,ly1,lcY>\DeclareUnicodeCharacter{00F8}{\o}
253 <all,t1,ly1>\DeclareUnicodeCharacter{00F9}{\@tabacckludge'u}
254 <all,t1,ly1>\DeclareUnicodeCharacter{00FA}{\@tabacckludge'u}
255 <all,t1,ly1>\DeclareUnicodeCharacter{00FB}{\^u}
256 <all,t1,ly1>\DeclareUnicodeCharacter{00FC}{\"u}
257 <all,t1,ly1>\DeclareUnicodeCharacter{00FD}{\@tabacckludge'y}
258 <all,t1,ly1>\DeclareUnicodeCharacter{00FE}{\th}
259 <all,t1,ly1>\DeclareUnicodeCharacter{00FF}{\"y}
260 <all,t1>\DeclareUnicodeCharacter{0102}{\u A}
261 <all,t1>\DeclareUnicodeCharacter{0103}{\u a}
262 <all,t1>\DeclareUnicodeCharacter{0104}{\k A}
263 <all,t1>\DeclareUnicodeCharacter{0105}{\k a}
264 <all,t1>\DeclareUnicodeCharacter{0106}{\@tabacckludge'C}
265 <all,t1>\DeclareUnicodeCharacter{0107}{\@tabacckludge'c}
266 <all,t1>\DeclareUnicodeCharacter{010C}{\v C}
267 <all,t1>\DeclareUnicodeCharacter{010D}{\v c}
268 <all,t1>\DeclareUnicodeCharacter{010E}{\v D}
269 <all,t1>\DeclareUnicodeCharacter{010F}{\v d}

```

```

270 <all,t1>\DeclareUnicodeCharacter{0110}{\DJ}
271 <all,t1>\DeclareUnicodeCharacter{0111}{\dj}
272 <all,t1>\DeclareUnicodeCharacter{0118}{\k E}
273 <all,t1>\DeclareUnicodeCharacter{0119}{\k e}
274 <all,t1>\DeclareUnicodeCharacter{011A}{\v E}
275 <all,t1>\DeclareUnicodeCharacter{011B}{\v e}
276 <all,t1>\DeclareUnicodeCharacter{011E}{\u G}
277 <all,t1>\DeclareUnicodeCharacter{011F}{\u g}
278 <all,t1>\DeclareUnicodeCharacter{0130}{\.I}
279 <all,t2c,t2b,t2a,t1,ot2,ot1,ly1,lcy>\DeclareUnicodeCharacter{0131}{\i}
280 <all,t1>\DeclareUnicodeCharacter{0132}{\IJ}
281 <all,t1>\DeclareUnicodeCharacter{0133}{\ij}
282 <all,t1>\DeclareUnicodeCharacter{0139}{\@tabacckludge'L}
283 <all,t1>\DeclareUnicodeCharacter{013A}{\@tabacckludge'l}
284 <all,t1>\DeclareUnicodeCharacter{013D}{\v L}
285 <all,t1>\DeclareUnicodeCharacter{013E}{\v l}
286 <all,t1,ot1,ly1>\DeclareUnicodeCharacter{0141}{\L}
287 <all,t1,ot1,ly1>\DeclareUnicodeCharacter{0142}{\l}
288 <all,t1>\DeclareUnicodeCharacter{0143}{\@tabacckludge'N}
289 <all,t1>\DeclareUnicodeCharacter{0144}{\@tabacckludge'n}
290 <all,t1>\DeclareUnicodeCharacter{0147}{\v N}
291 <all,t1>\DeclareUnicodeCharacter{0148}{\v n}
292 <all,t1>\DeclareUnicodeCharacter{014A}{\NG}
293 <all,t1>\DeclareUnicodeCharacter{014B}{\ng}
294 <all,t1>\DeclareUnicodeCharacter{0150}{\H O}
295 <all,t1>\DeclareUnicodeCharacter{0151}{\H o}
296 <all,t1,ot1,ly1,lcy>\DeclareUnicodeCharacter{0152}{\OE}
297 <all,t1,ot1,ly1,lcy>\DeclareUnicodeCharacter{0153}{\oe}
298 <all,t1>\DeclareUnicodeCharacter{0154}{\@tabacckludge'R}
299 <all,t1>\DeclareUnicodeCharacter{0155}{\@tabacckludge'r}
300 <all,t1>\DeclareUnicodeCharacter{0158}{\v R}
301 <all,t1>\DeclareUnicodeCharacter{0159}{\v r}
302 <all,t1>\DeclareUnicodeCharacter{015A}{\@tabacckludge'S}
303 <all,t1>\DeclareUnicodeCharacter{015B}{\@tabacckludge's}
304 <all,t1>\DeclareUnicodeCharacter{015E}{\c S}
305 <all,t1>\DeclareUnicodeCharacter{015F}{\c s}
306 <all,t1,ly1>\DeclareUnicodeCharacter{0160}{\v S}
307 <all,t1,ly1>\DeclareUnicodeCharacter{0161}{\v s}
308 <all,t1>\DeclareUnicodeCharacter{0162}{\c T}
309 <all,t1>\DeclareUnicodeCharacter{0163}{\c t}
310 <all,t1>\DeclareUnicodeCharacter{0164}{\v T}
311 <all,t1>\DeclareUnicodeCharacter{0165}{\v t}
312 <all,t1>\DeclareUnicodeCharacter{016E}{\r U}
313 <all,t1>\DeclareUnicodeCharacter{016F}{\r u}
314 <all,t1>\DeclareUnicodeCharacter{0170}{\H U}
315 <all,t1>\DeclareUnicodeCharacter{0171}{\H u}
316 <all,t1,ly1>\DeclareUnicodeCharacter{0178}{\ "Y}
317 <all,t1>\DeclareUnicodeCharacter{0179}{\@tabacckludge'Z}
318 <all,t1>\DeclareUnicodeCharacter{017A}{\@tabacckludge'z}
319 <all,t1>\DeclareUnicodeCharacter{017B}{\ .Z}
320 <all,t1>\DeclareUnicodeCharacter{017C}{\ .z}
321 <all,t1,ly1>\DeclareUnicodeCharacter{017D}{\v Z}
322 <all,t1,ly1>\DeclareUnicodeCharacter{017E}{\v z}
323 <all,ts1,ly1>\DeclareUnicodeCharacter{0192}{\textflorin}

```

```

324 <all, ly1, utf8>\DeclareUnicodeCharacter{02C6}{\textasciicircum}
325 <all, ts1>\DeclareUnicodeCharacter{02C7}{\textasciicaron}
326 <all, ly1, utf8>\DeclareUnicodeCharacter{02DC}{\textasciitilde}
327 <all, ts1>\DeclareUnicodeCharacter{02D8}{\textasciibreve}
328 <all, ts1>\DeclareUnicodeCharacter{02DD}{\textacutedbl}

```

The Cyrillic code points have been recently checked (2007) and extended and corrected by Matthias Noe (a9931078@unet.univie.ac.at) — thanks.

```

329 <*all, x2, t2c, t2b, t2a, ot2, lcy>
330 \DeclareUnicodeCharacter{0400}{\@tabacckludge'\CYRE}
331 </all, x2, t2c, t2b, t2a, ot2, lcy>
332 <all, x2, t2c, t2b, t2a, ot2, lcy>\DeclareUnicodeCharacter{0401}{\CYRYO}
333 <all, x2, t2a, ot2>\DeclareUnicodeCharacter{0402}{\CYRDJE}
334 <*all, x2, t2c, t2b, t2a, ot2, lcy>
335 \DeclareUnicodeCharacter{0403}{\@tabacckludge'\CYRG}
336 </all, x2, t2c, t2b, t2a, ot2, lcy>
337 <all, x2, t2a, ot2, lcy>\DeclareUnicodeCharacter{0404}{\CYRIE}
338 <all, x2, t2c, t2b, t2a, ot2>\DeclareUnicodeCharacter{0405}{\CYRDZE}
339 <all, x2, t2c, t2b, t2a, ot2, lcy>\DeclareUnicodeCharacter{0406}{\CYRII}
340 <all, x2, t2a, lcy>\DeclareUnicodeCharacter{0407}{\CYRYI}
341 <all, x2, t2c, t2b, t2a, ot2>\DeclareUnicodeCharacter{0408}{\CYRJE}
342 <all, x2, t2b, t2a, ot2>\DeclareUnicodeCharacter{0409}{\CYRLJE}
343 <all, x2, t2b, t2a, ot2>\DeclareUnicodeCharacter{040A}{\CYRNJE}
344 <all, x2, t2a, ot2>\DeclareUnicodeCharacter{040B}{\CYRTSHE}
345 <*all, x2, t2c, t2b, t2a, ot2, lcy>
346 \DeclareUnicodeCharacter{040C}{\@tabacckludge'\CYRK}
347 \DeclareUnicodeCharacter{040D}{\@tabacckludge'\CYRI}
348 </all, x2, t2c, t2b, t2a, ot2, lcy>
349 <all, x2, t2b, t2a, lcy>\DeclareUnicodeCharacter{040E}{\CYRUSHRT}
350 <all, x2, t2c, t2a, ot2>\DeclareUnicodeCharacter{040F}{\CYRDZHE}
351 <*all, x2, t2c, t2b, t2a, ot2, lcy>
352 \DeclareUnicodeCharacter{0410}{\CYRA}
353 \DeclareUnicodeCharacter{0411}{\CYRB}
354 \DeclareUnicodeCharacter{0412}{\CYRV}
355 \DeclareUnicodeCharacter{0413}{\CYRG}
356 \DeclareUnicodeCharacter{0414}{\CYRD}
357 \DeclareUnicodeCharacter{0415}{\CYRE}
358 \DeclareUnicodeCharacter{0416}{\CYRZH}
359 \DeclareUnicodeCharacter{0417}{\CYRZ}
360 \DeclareUnicodeCharacter{0418}{\CYRI}
361 \DeclareUnicodeCharacter{0419}{\CYRISHRT}
362 \DeclareUnicodeCharacter{041A}{\CYRK}
363 \DeclareUnicodeCharacter{041B}{\CYRL}
364 \DeclareUnicodeCharacter{041C}{\CYRM}
365 \DeclareUnicodeCharacter{041D}{\CYRN}
366 \DeclareUnicodeCharacter{041E}{\CYRO}
367 \DeclareUnicodeCharacter{041F}{\CYRP}
368 \DeclareUnicodeCharacter{0420}{\CYRR}
369 \DeclareUnicodeCharacter{0421}{\CYRS}
370 \DeclareUnicodeCharacter{0422}{\CYRT}
371 \DeclareUnicodeCharacter{0423}{\CYRU}
372 \DeclareUnicodeCharacter{0424}{\CYRF}
373 \DeclareUnicodeCharacter{0425}{\CYRH}
374 \DeclareUnicodeCharacter{0426}{\CYRC}

```

```

375 \DeclareUnicodeCharacter{0427}{\CYRCH}
376 \DeclareUnicodeCharacter{0428}{\CYRSH}
377 \DeclareUnicodeCharacter{0429}{\CYRSHCH}
378 \DeclareUnicodeCharacter{042A}{\CYRHRDSN}
379 \DeclareUnicodeCharacter{042B}{\CYRERY}
380 \DeclareUnicodeCharacter{042C}{\CYRSFTSN}
381 \DeclareUnicodeCharacter{042D}{\CYREREV}
382 \DeclareUnicodeCharacter{042E}{\CYRYU}
383 \DeclareUnicodeCharacter{042F}{\CYRYA}
384 \DeclareUnicodeCharacter{0430}{\cyra}
385 \DeclareUnicodeCharacter{0431}{\cyrb}
386 \DeclareUnicodeCharacter{0432}{\cyrv}
387 \DeclareUnicodeCharacter{0433}{\cyrg}
388 \DeclareUnicodeCharacter{0434}{\cyrd}
389 \DeclareUnicodeCharacter{0435}{\cyre}
390 \DeclareUnicodeCharacter{0436}{\cyrzh}
391 \DeclareUnicodeCharacter{0437}{\cyrz}
392 \DeclareUnicodeCharacter{0438}{\cyri}
393 \DeclareUnicodeCharacter{0439}{\cyrishrt}
394 \DeclareUnicodeCharacter{043A}{\cyrk}
395 \DeclareUnicodeCharacter{043B}{\cyr1}
396 \DeclareUnicodeCharacter{043C}{\cyrn}
397 \DeclareUnicodeCharacter{043D}{\cyrn}
398 \DeclareUnicodeCharacter{043E}{\cyro}
399 \DeclareUnicodeCharacter{043F}{\cyrp}
400 \DeclareUnicodeCharacter{0440}{\cyrr}
401 \DeclareUnicodeCharacter{0441}{\cyrs}
402 \DeclareUnicodeCharacter{0442}{\cyrst}
403 \DeclareUnicodeCharacter{0443}{\cyru}
404 \DeclareUnicodeCharacter{0444}{\cyrf}
405 \DeclareUnicodeCharacter{0445}{\cyrh}
406 \DeclareUnicodeCharacter{0446}{\cyrc}
407 \DeclareUnicodeCharacter{0447}{\cyrch}
408 \DeclareUnicodeCharacter{0448}{\cyrsh}
409 \DeclareUnicodeCharacter{0449}{\cyrshch}
410 \DeclareUnicodeCharacter{044A}{\cyrhrdsn}
411 \DeclareUnicodeCharacter{044B}{\cyrery}
412 \DeclareUnicodeCharacter{044C}{\cyrstsn}
413 \DeclareUnicodeCharacter{044D}{\cyrerev}
414 \DeclareUnicodeCharacter{044E}{\cyryu}
415 \DeclareUnicodeCharacter{044F}{\cyrya}
416 \DeclareUnicodeCharacter{0450}{\@tabacckludge'\cyre}
417 \DeclareUnicodeCharacter{0451}{\cyryo}
418 </all,x2,t2c,t2b,t2a,ot2,lcy>
419 <all,x2,t2a,ot2>\DeclareUnicodeCharacter{0452}{\cyrdje}
420 <*all,x2,t2c,t2b,t2a,ot2,lcy>
421 \DeclareUnicodeCharacter{0453}{\@tabacckludge'\cyrg}
422 </all,x2,t2c,t2b,t2a,ot2,lcy>
423 <all,x2,t2a,ot2,lcy>\DeclareUnicodeCharacter{0454}{\cyrie}
424 <all,x2,t2c,t2b,t2a,ot2>\DeclareUnicodeCharacter{0455}{\cyrdze}
425 <all,x2,t2c,t2b,t2a,ot2,lcy>\DeclareUnicodeCharacter{0456}{\cyr11}
426 <all,x2,t2a,lcy>\DeclareUnicodeCharacter{0457}{\cyr1i}
427 <all,x2,t2c,t2b,t2a,ot2>\DeclareUnicodeCharacter{0458}{\cyr1je}
428 <all,x2,t2b,t2a,ot2>\DeclareUnicodeCharacter{0459}{\cyr1je}

```

```

429 <all, x2, t2b, t2a, ot2>\DeclareUnicodeCharacter{045A}{\cyrnje}
430 <all, x2, t2a, ot2>\DeclareUnicodeCharacter{045B}{\cyrtshe}
431 <*all, x2, t2c, t2b, t2a, ot2, lcy>
432 \DeclareUnicodeCharacter{045C}{\@tabacckludge'\cyrk}
433 \DeclareUnicodeCharacter{045D}{\@tabacckludge'\cyri}
434 </all, x2, t2c, t2b, t2a, ot2, lcy>
435 <all, x2, t2b, t2a, lcy>\DeclareUnicodeCharacter{045E}{\cyrushrt}
436 <all, x2, t2c, t2a, ot2>\DeclareUnicodeCharacter{045F}{\cyrdzhe}
437 <all, x2, ot2>\DeclareUnicodeCharacter{0462}{\CYRYAT}
438 <all, x2, ot2>\DeclareUnicodeCharacter{0463}{\cyryat}
439 <all, x2>\DeclareUnicodeCharacter{046A}{\CYRBYUS}
440 <all, x2>\DeclareUnicodeCharacter{046B}{\cyrbyus}

```

The next two declarations are questionable, the encoding definition should probably contain \CYROTLD and \cyrotld. Or alternatively, if the characters in the X2 encodings are really meant to represent the historical characters in Ux0472 and Ux0473 (they look like them) then they would need to change instead.

However, their looks are probably a font designers decision and the next two mappings are wrong or rather the names in OT2 should change for consistency.

On the other hand the names \CYROTLD are somewhat questionabled as the Unicode standard only describes “Cyrillic barred O” while TLD refers to a tilde (which is more less what the “Cyrillic FITA looks according to the Unicode book).

```

441 <all, ot2>\DeclareUnicodeCharacter{0472}{\CYRFITA}
442 <all, ot2>\DeclareUnicodeCharacter{0473}{\cyrfita}
443 <all, x2, ot2>\DeclareUnicodeCharacter{0474}{\CYRIZH}
444 <all, x2, ot2>\DeclareUnicodeCharacter{0475}{\cyrizh}

```

While the double grave accent seems to exist in X2, T2A, T2B and T2C encoding, the letter izhitsa exists only in X2 and OT2. Therefore, izhitsa with double grave seems to be possible only using X2.

```

445 <all, x2>\DeclareUnicodeCharacter{0476}{\C\CYRIZH}
446 <all, x2>\DeclareUnicodeCharacter{0477}{\C\cyrizh}
447 <all, t2c>\DeclareUnicodeCharacter{048C}{\CYRSEMISFTSN}
448 <all, t2c>\DeclareUnicodeCharacter{048D}{\cyrsemisftsn}
449 <all, t2c>\DeclareUnicodeCharacter{048E}{\CYRRTICK}
450 <all, t2c>\DeclareUnicodeCharacter{048F}{\cyrrtick}
451 <all, x2, t2a, lcy>\DeclareUnicodeCharacter{0490}{\CYRGUP}
452 <all, x2, t2a, lcy>\DeclareUnicodeCharacter{0491}{\cyrgup}
453 <all, x2, t2b, t2a>\DeclareUnicodeCharacter{0492}{\CYRGHCRS}
454 <all, x2, t2b, t2a>\DeclareUnicodeCharacter{0493}{\cyrghcrs}
455 <all, x2, t2c, t2b>\DeclareUnicodeCharacter{0494}{\CYRGHK}
456 <all, x2, t2c, t2b>\DeclareUnicodeCharacter{0495}{\cyrghk}
457 <all, x2, t2b, t2a>\DeclareUnicodeCharacter{0496}{\CYRZHDSC}
458 <all, x2, t2b, t2a>\DeclareUnicodeCharacter{0497}{\cyrzhdsc}
459 <all, x2, t2a>\DeclareUnicodeCharacter{0498}{\CYRZDSC}
460 <all, x2, t2a>\DeclareUnicodeCharacter{0499}{\cyrzdsc}
461 <all, x2, t2c, t2b, t2a>\DeclareUnicodeCharacter{049A}{\CYRKDSC}
462 <all, x2, t2c, t2b, t2a>\DeclareUnicodeCharacter{049B}{\cyrkdsc}
463 <all, x2, t2a>\DeclareUnicodeCharacter{049C}{\CYRKVCRS}
464 <all, x2, t2a>\DeclareUnicodeCharacter{049D}{\cyrkvcrs}
465 <all, x2, t2c>\DeclareUnicodeCharacter{049E}{\CYRKHCRS}
466 <all, x2, t2c>\DeclareUnicodeCharacter{049F}{\cyrkhcrs}
467 <all, x2, t2a>\DeclareUnicodeCharacter{04A0}{\CYRKBEAK}

```



```

468 <all, x2, t2a>\DeclareUnicodeCharacter{04A1}{\cyrkbeak}
469 <all, x2, t2c, t2b, t2a>\DeclareUnicodeCharacter{04A2}{\CYRNDSC}
470 <all, x2, t2c, t2b, t2a>\DeclareUnicodeCharacter{04A3}{\cyrndsc}
471 <all, x2, t2b, t2a>\DeclareUnicodeCharacter{04A4}{\CYRNG}
472 <all, x2, t2b, t2a>\DeclareUnicodeCharacter{04A5}{\cyrng}
473 <all, x2, t2c>\DeclareUnicodeCharacter{04A6}{\CYRPHK}
474 <all, x2, t2c>\DeclareUnicodeCharacter{04A7}{\cyrphk}
475 <all, x2, t2c>\DeclareUnicodeCharacter{04A8}{\CYRABHHA}
476 <all, x2, t2c>\DeclareUnicodeCharacter{04A9}{\cyrabhha}
477 <all, x2, t2a>\DeclareUnicodeCharacter{04AA}{\CYRSDSC}
478 <all, x2, t2a>\DeclareUnicodeCharacter{04AB}{\cyrsdsc}
479 <all, x2, t2c>\DeclareUnicodeCharacter{04AC}{\CYRTDSC}
480 <all, x2, t2c>\DeclareUnicodeCharacter{04AD}{\cyrt dsc}
481 <all, x2, t2b, t2a>\DeclareUnicodeCharacter{04AE}{\CYRY}
482 <all, x2, t2b, t2a>\DeclareUnicodeCharacter{04AF}{\cyry}
483 <all, x2, t2a>\DeclareUnicodeCharacter{04B0}{\CYRYHCRS}
484 <all, x2, t2a>\DeclareUnicodeCharacter{04B1}{\cyryhcrs}
485 <all, x2, t2c, t2b, t2a>\DeclareUnicodeCharacter{04B2}{\CYRHDSC}
486 <all, x2, t2c, t2b, t2a>\DeclareUnicodeCharacter{04B3}{\cyrhdsc}
487 <all, x2, t2c>\DeclareUnicodeCharacter{04B4}{\CYRTETSE}
488 <all, x2, t2c>\DeclareUnicodeCharacter{04B5}{\cyrtetse}
489 <all, x2, t2c, t2b, t2a>\DeclareUnicodeCharacter{04B6}{\CYRCHRDSC}
490 <all, x2, t2c, t2b, t2a>\DeclareUnicodeCharacter{04B7}{\cyrchrdsc}
491 <all, x2, t2a>\DeclareUnicodeCharacter{04B8}{\CYRCHVCRS}
492 <all, x2, t2a>\DeclareUnicodeCharacter{04B9}{\cyrchvcrs}
493 <all, x2, t2c, t2b, t2a>\DeclareUnicodeCharacter{04BA}{\CYRSHHA}
494 <all, x2, t2c, t2b, t2a>\DeclareUnicodeCharacter{04BB}{\cyrshha}
495 <all, x2, t2c>\DeclareUnicodeCharacter{04BC}{\CYRABHCH}
496 <all, x2, t2c>\DeclareUnicodeCharacter{04BD}{\cyrabhch}
497 <all, x2, t2c>\DeclareUnicodeCharacter{04BE}{\CYRABHCHDSC}
498 <all, x2, t2c>\DeclareUnicodeCharacter{04BF}{\cyrabhchdsc}

```

The character \CYRpalochka is not defined by OT2 and LCY. However it is looking identical to \CYRII and the Unicode standard explicitly refers to that (and to Latin I). So perhaps those encodings could get an alias? On the other hand, why are there two distinct slots in the T2 encodings even though they are so pressed for space? Perhaps they don't always look alike.

```

499 <all, x2, t2c, t2b, t2a>\DeclareUnicodeCharacter{04C0}{\CYRpalochka}
500 <all, x2, t2c, t2b, t2a, ot2, lcy>\DeclareUnicodeCharacter{04C1}{\U\CYRZH}
501 <all, x2, t2c, t2b, t2a, ot2, lcy>\DeclareUnicodeCharacter{04C2}{\U\cyrzh}
502 <all, x2, t2b>\DeclareUnicodeCharacter{04C3}{\CYRKHK}
503 <all, x2, t2b>\DeclareUnicodeCharacter{04C4}{\cyrkhk}

```

According to the Unicode standard Ux04C5 should be an L with “tail” not with descender (which also exists as Ux04A2) but it looks as if the char names do not make this distinction). Should they?

```

504 <all, x2, t2c, t2b>\DeclareUnicodeCharacter{04C5}{\CYRLDSC}
505 <all, x2, t2c, t2b>\DeclareUnicodeCharacter{04C6}{\cyrl dsc}
506 <all, x2, t2c, t2b>\DeclareUnicodeCharacter{04C7}{\CYRNHK}
507 <all, x2, t2c, t2b>\DeclareUnicodeCharacter{04C8}{\cyrnhk}
508 <all, x2, t2b>\DeclareUnicodeCharacter{04CB}{\CYRCHLDSC}
509 <all, x2, t2b>\DeclareUnicodeCharacter{04CC}{\cyrchldsc}

```

According to the Unicode standard Ux04CD should be an M with “tail” not with descender. However this time there is no M with descender in the Unicode standard.

```

510 (all, x2, t2c)\DeclareUnicodeCharacter{04CD}{\CYRMDSC}
511 (all, x2, t2c)\DeclareUnicodeCharacter{04CE}{\cyrmdsc}

512 (all, x2, t2c, t2b, t2a, ot2, lcy)\DeclareUnicodeCharacter{04D0}{\U\CYRA}
513 (all, x2, t2c, t2b, t2a, ot2, lcy)\DeclareUnicodeCharacter{04D1}{\U\cyra}
514 (all, x2, t2c, t2b, t2a, ot2, lcy)\DeclareUnicodeCharacter{04D2}{\\"CYRA}
515 (all, x2, t2c, t2b, t2a, ot2, lcy)\DeclareUnicodeCharacter{04D3}{\\"cyra}
516 (all, x2, t2a)\DeclareUnicodeCharacter{04D4}{\CYRAE}
517 (all, x2, t2a)\DeclareUnicodeCharacter{04D5}{\cyrae}
518 (all, x2, t2c, t2b, t2a, ot2, lcy)\DeclareUnicodeCharacter{04D6}{\U\CYRE}
519 (all, x2, t2c, t2b, t2a, ot2, lcy)\DeclareUnicodeCharacter{04D7}{\U\cyre}
520 (all, x2, t2c, t2b, t2a)\DeclareUnicodeCharacter{04D8}{\CYRSCHWA}
521 (all, x2, t2c, t2b, t2a)\DeclareUnicodeCharacter{04D9}{\cyrschwa}
522 (all, x2, t2c, t2b, t2a)\DeclareUnicodeCharacter{04DA}{\\"CYRSCHWA}
523 (all, x2, t2c, t2b, t2a)\DeclareUnicodeCharacter{04DB}{\\"cyrschwa}
524 (all, x2, t2c, t2b, t2a, ot2, lcy)\DeclareUnicodeCharacter{04DC}{\\"CYRZH}
525 (all, x2, t2c, t2b, t2a, ot2, lcy)\DeclareUnicodeCharacter{04DD}{\\"cyrz}
526 (all, x2, t2c, t2b, t2a, ot2, lcy)\DeclareUnicodeCharacter{04DE}{\\"CYRZ}
527 (all, x2, t2c, t2b, t2a, ot2, lcy)\DeclareUnicodeCharacter{04DF}{\\"cyrz}
528 (all, x2, t2c, t2b)\DeclareUnicodeCharacter{04E0}{\CYRABHDZE}
529 (all, x2, t2c, t2b)\DeclareUnicodeCharacter{04E1}{\cyrabhdze}
530 (all, x2, t2c, t2b, t2a, ot2, lcy)\DeclareUnicodeCharacter{04E2}{\CYRI}
531 (all, x2, t2c, t2b, t2a, ot2, lcy)\DeclareUnicodeCharacter{04E3}{\cyri}
532 (all, x2, t2c, t2b, t2a, ot2, lcy)\DeclareUnicodeCharacter{04E4}{\\"CYRI}
533 (all, x2, t2c, t2b, t2a, ot2, lcy)\DeclareUnicodeCharacter{04E5}{\\"cyri}
534 (all, x2, t2c, t2b, t2a, ot2, lcy)\DeclareUnicodeCharacter{04E6}{\\"CYRO}
535 (all, x2, t2c, t2b, t2a, ot2, lcy)\DeclareUnicodeCharacter{04E7}{\\"cyro}
536 (all, x2, t2c, t2b, t2a)\DeclareUnicodeCharacter{04E8}{\CYROTLD}
537 (all, x2, t2c, t2b, t2a)\DeclareUnicodeCharacter{04E9}{\cyrotld}
538 (all, x2, t2c, t2b, t2a, ot2, lcy)\DeclareUnicodeCharacter{04EC}{\\"CYREREV}
539 (all, x2, t2c, t2b, t2a, ot2, lcy)\DeclareUnicodeCharacter{04ED}{\\"cyrerev}
540 (all, x2, t2c, t2b, t2a, ot2, lcy)\DeclareUnicodeCharacter{04EE}{\CYRU}
541 (all, x2, t2c, t2b, t2a, ot2, lcy)\DeclareUnicodeCharacter{04EF}{\cyru}
542 (all, x2, t2c, t2b, t2a, ot2, lcy)\DeclareUnicodeCharacter{04F0}{\\"CYRU}
543 (all, x2, t2c, t2b, t2a, ot2, lcy)\DeclareUnicodeCharacter{04F1}{\\"cyru}
544 (all, x2, t2c, t2b, t2a, ot2, lcy)\DeclareUnicodeCharacter{04F2}{\H\CYRU}
545 (all, x2, t2c, t2b, t2a, ot2, lcy)\DeclareUnicodeCharacter{04F3}{\H\cyru}
546 (all, x2, t2c, t2b, t2a, ot2, lcy)\DeclareUnicodeCharacter{04F4}{\\"CYRCH}
547 (all, x2, t2c, t2b, t2a, ot2, lcy)\DeclareUnicodeCharacter{04F5}{\\"cyrch}
548 (all, x2, t2b)\DeclareUnicodeCharacter{04F6}{\CYRGDSC}
549 (all, x2, t2b)\DeclareUnicodeCharacter{04F7}{\cyrgdsc}
550 (all, x2, t2c, t2b, t2a, ot2, lcy)\DeclareUnicodeCharacter{04F8}{\\"CYRERY}
551 (all, x2, t2c, t2b, t2a, ot2, lcy)\DeclareUnicodeCharacter{04F9}{\\"cyrery}
552 (all, t2b)\DeclareUnicodeCharacter{04FA}{\CYRGDSCHCRS}
553 (all, t2b)\DeclareUnicodeCharacter{04FB}{\cyrgdschcrs}
554 (all, x2, t2b)\DeclareUnicodeCharacter{04FC}{\CYRHHK}
555 (all, x2, t2b)\DeclareUnicodeCharacter{04FD}{\cyrrhk}
556 (all, t2b)\DeclareUnicodeCharacter{04FE}{\CYRHHCRS}
557 (all, t2b)\DeclareUnicodeCharacter{04FF}{\cyrrhcrs}
558 (all, ts1)\DeclareUnicodeCharacter{0E3F}{\textbaht}
559 (all, x2, t2c, t2b, t2a, t1, utf8)\DeclareUnicodeCharacter{200C}{\textcompwordmark}

```

```

560 (*all,x2,t2c,t2b,t2a,t1,ot2,ot1,ly1,lcy)
561 \DeclareUnicodeCharacter{2013}{\textendash}
562 \DeclareUnicodeCharacter{2014}{\textemdash}
563 (/all,x2,t2c,t2b,t2a,t1,ot2,ot1,ly1,lcy)
564 (all,ts1)\DeclareUnicodeCharacter{2016}{\textbardbl}
565 (*all,x2,t2c,t2b,t2a,t1,ot2,ot1,lcy)
566 \DeclareUnicodeCharacter{2018}{\textquoteleft}
567 \DeclareUnicodeCharacter{2019}{\textquoteright}
568 (/all,x2,t2c,t2b,t2a,t1,ot2,ot1,lcy)
569 (all,t1)\DeclareUnicodeCharacter{201A}{\quotesinglbase}
570 (*all,x2,t2c,t2b,t2a,t1,ot2,ot1,ly1,lcy)
571 \DeclareUnicodeCharacter{201C}{\textquotedblleft}
572 \DeclareUnicodeCharacter{201D}{\textquotedblright}
573 (/all,x2,t2c,t2b,t2a,t1,ot2,ot1,ly1,lcy)
574 (all,x2,t2c,t2b,t2a,t1,lcy)\DeclareUnicodeCharacter{201E}{\quotedblbase}
575 (all,ts1,oms,ly1)\DeclareUnicodeCharacter{2020}{\textdagger}
576 (all,ts1,oms,ly1)\DeclareUnicodeCharacter{2021}{\textdaggerdbl}
577 (all,ts1,oms,ly1)\DeclareUnicodeCharacter{2022}{\textbullet}
578 (all,ly1,utf8)\DeclareUnicodeCharacter{2026}{\textellipsis}
579 (*all,x2,ts1,t2c,t2b,t2a,t1,ly1)
580 \DeclareUnicodeCharacter{2030}{\textperthousand}
581 (/all,x2,ts1,t2c,t2b,t2a,t1,ly1)
582 (*all,x2,ts1,t2c,t2b,t2a,t1)
583 \DeclareUnicodeCharacter{2031}{\textpertenthousand}
584 (/all,x2,ts1,t2c,t2b,t2a,t1)
585 (all,t1,ly1)\DeclareUnicodeCharacter{2039}{\guilsinglleft}
586 (all,t1,ly1)\DeclareUnicodeCharacter{203A}{\guilsinglright}
587 (all,ts1)\DeclareUnicodeCharacter{203B}{\textreferencemark}
588 (all,ts1)\DeclareUnicodeCharacter{203D}{\textinterrobang}
589 (all,ts1)\DeclareUnicodeCharacter{2044}{\textfractionsolidus}
590 (all,ts1)\DeclareUnicodeCharacter{204E}{\textasteriskcentered}
591 (all,ts1)\DeclareUnicodeCharacter{2052}{\textdiscount}
592 (all,ts1)\DeclareUnicodeCharacter{20A1}{\textcolonmonetary}
593 (all,ts1)\DeclareUnicodeCharacter{20A4}{\textlira}
594 (all,ts1)\DeclareUnicodeCharacter{20A6}{\textnaira}
595 (all,ts1)\DeclareUnicodeCharacter{20A9}{\textwon}
596 (all,ts1)\DeclareUnicodeCharacter{20AB}{\textdong}
597 (all,ts1)\DeclareUnicodeCharacter{20AC}{\texteuro}
598 (all,ts1)\DeclareUnicodeCharacter{20B1}{\textpeso}
599 (all,ts1)\DeclareUnicodeCharacter{2103}{\textcelsius}
600 (all,x2,ts1,t2c,t2b,t2a,ot2,lcy)\DeclareUnicodeCharacter{2116}{\textnumero}
601 (all,ts1)\DeclareUnicodeCharacter{2117}{\textcircledP}
602 (all,ts1)\DeclareUnicodeCharacter{211E}{\textrecipe}
603 (all,ts1)\DeclareUnicodeCharacter{2120}{\textservicemark}
604 (all,ts1,ly1,utf8)\DeclareUnicodeCharacter{2122}{\texttrademark}
605 (all,ts1)\DeclareUnicodeCharacter{2126}{\textohm}
606 (all,ts1)\DeclareUnicodeCharacter{2127}{\textmho}
607 (all,ts1)\DeclareUnicodeCharacter{212E}{\textestimated}
608 (all,ts1)\DeclareUnicodeCharacter{2190}{\textleftarrow}
609 (all,ts1)\DeclareUnicodeCharacter{2191}{\textuparrow}
610 (all,ts1)\DeclareUnicodeCharacter{2192}{\textrightarrow}
611 (all,ts1)\DeclareUnicodeCharacter{2193}{\textdownarrow}
612 (all,x2,ts1,t2c,t2b,t2a)\DeclareUnicodeCharacter{2329}{\texttriangle}
613 (all,x2,ts1,t2c,t2b,t2a)\DeclareUnicodeCharacter{232A}{\texttriangle}

```

```

614 <all,ts1>\DeclareUnicodeCharacter{2422}{\textblank}
615 <all,x2,t2c,t2b,t2a,t1,utf8>\DeclareUnicodeCharacter{2423}{\textvisiblespace}
616 <all,ts1>\DeclareUnicodeCharacter{25E6}{\textopenbullet}
617 <all,ts1>\DeclareUnicodeCharacter{25EF}{\textbigcirc}
618 <all,ts1>\DeclareUnicodeCharacter{266A}{\textmusicalnote}

```

3.3 Notes

The following inputs are inconsistent with the 8-bit inputenc files since they will always only produce the ‘text character’. This is an area where inputenc is notoriously confused.

```

%<all,ts1,t1,ot1,ly1>\DeclareUnicodeCharacter{00A3}{\textsterling}
%<*all,x2,ts1,t2c,t2b,t2a,oms,ly1>
\DeclareUnicodeCharacter{00A7}{\textsection}
%</all,x2,ts1,t2c,t2b,t2a,oms,ly1>
%<all,ts1,utf8>\DeclareUnicodeCharacter{00A9}{\textcopyright}
%<all,ts1>\DeclareUnicodeCharacter{00B1}{\textpm}
%<all,ts1,oms,ly1>\DeclareUnicodeCharacter{00B6}{\textparagraph}
%<all,ts1,oms,ly1>\DeclareUnicodeCharacter{2020}{\textdagger}
%<all,ts1,oms,ly1>\DeclareUnicodeCharacter{2021}{\textdaggerdbl}
%<all,ly1,utf8>\DeclareUnicodeCharacter{2026}{\textellipsis}

```

The following definitions are in an encoding file but have no direct equivalent in Unicode, or they simply do not make sense in that context (or we have not yet found anything or ...:-). For example, the non-combining accent characters are certainly available somewhere but these are not equivalent to a \TeX accent command.

```

\DeclareTextSymbol{\j}{OT1}{17}
\DeclareTextSymbol{\SS}{T1}{223}
\DeclareTextSymbol{\textcompwordmark}{T1}{23}

\DeclareTextAccent{"}{OT1}{127}
\DeclareTextAccent{\'}{OT1}{19}
\DeclareTextAccent{\.}{OT1}{95}
\DeclareTextAccent{\=}{OT1}{22}
\DeclareTextAccent{\H}{OT1}{125}
\DeclareTextAccent{\^}{OT1}{94}
\DeclareTextAccent{\'}{OT1}{18}
\DeclareTextAccent{\r}{OT1}{23}
\DeclareTextAccent{\u}{OT1}{21}
\DeclareTextAccent{\v}{OT1}{20}
\DeclareTextAccent{\~}{OT1}{126}
\DeclareTextCommand{\b}{OT1}[1]
\DeclareTextCommand{\c}{OT1}[1]
\DeclareTextCommand{\d}{OT1}[1]
\DeclareTextCommand{\k}{T1}[1]

```

3.4 Mappings for OT1 glyphs

This is even more incomplete as again it covers only the single glyphs from OT1 plus some that have been explicitly defined for this encoding. Everything that is

provided in T1, and that could be provided as composite glyphs via OT1, could and probably should be set up as well. Which leaves the many things that are not provided in T1 but can be provided in OT1 (and in T1) by composite glyphs.

Stuff not mapped (note that `\j` (*j*) is not equivalent to any Unicode character):

```
\DeclareTextSymbol{\j}{OT1}{17}
\DeclareTextAccent{"}{OT1}{127}
\DeclareTextAccent{'}{OT1}{19}
\DeclareTextAccent{.}{OT1}{95}
\DeclareTextAccent{=}{OT1}{22}
\DeclareTextAccent{^}{OT1}{94}
\DeclareTextAccent{'}{OT1}{18}
\DeclareTextAccent{~}{OT1}{126}
\DeclareTextAccent{\H}{OT1}{125}
\DeclareTextAccent{\u}{OT1}{21}
\DeclareTextAccent{\v}{OT1}{20}
\DeclareTextAccent{\r}{OT1}{23}
\DeclareTextCommand{\b}{OT1}[1]
\DeclareTextCommand{\c}{OT1}[1]
\DeclareTextCommand{\d}{OT1}[1]
```

3.5 Mappings for OMS glyphs

Characters like `\textbackslash` are not mapped as they are (primarily) only in the lower 127 and the code here only sets up mappings for UTF-8 characters that are at least 2 octets long.

```
\DeclareTextSymbol{\textbackslash}{OMS}{110}      % "6E
\DeclareTextSymbol{\textbar}{OMS}{106}            % "6A
\DeclareTextSymbol{\textbraceleft}{OMS}{102}      % "66
\DeclareTextSymbol{\textbraceright}{OMS}{103}     % "67
```

But the following (and some others) might actually lurk in Unicode somewhere...

```
\DeclareTextSymbol{\textasteriskcentered}{OMS}{3} % "03
\DeclareTextCommand{\textcircled}{OMS}
```

3.6 Mappings for TS1 glyphs

Exercise for somebody else.

3.7 Mappings for latex.ltx glyphs

There is also a collection of characters already set up in the kernel, one way or the other. Since these do not clearly relate to any particular font encoding they are mapped when the utf8 support is first set up.

Also there are a number of `\providecommands` in the various input encoding files which may or may not go into this part.

```
619 <*\utf8>
620 % This space is intentionally empty ...
621 </\utf8>
```

4 A test document

Here is a very small test document which may or may not survive if the current document is transferred from one place to the other.

```
622 (*test)
623 \documentclass{article}
624
625 \usepackage[latin1,utf8]{inputenc}
626 \usepackage[T1]{fontenc}
627 \usepackage{trace}
628
629 \scrollmode % to run past the error below
630
631 \begin{document}
632
633 German umlauts in UTF-8: ^^c3^^a4^^c3^^b6^^c3^^bc  %%% äöü
634
635 \inputencoding{latin1} % switch to latin1
636
637 German umlauts in UTF-8 but read by latin1 (and will produce one
638 error since \verb=\textcurrency= is not provided):
639 ^^c3^^a4^^c3^^b6^^c3^^bc
640
641 \inputencoding{utf8} % switch back to utf8
642
643 German umlauts in UTF-8: ^^c3^^a4^^c3^^b6^^c3^^bc
644
645
646 Some codes that should produce errors as nothing is set up
647 for them: ^^c3F ^^e1^^a4^^b6
648
649 And some that are not legal utf8 sequences: ^^c3X ^^e1XY
650
651 \showoutput
652 \tracingstats=2
653 \stop
654 /\test)
```