

Providing some UTF-8 support via `inputenc`

David Carlisle Frank Mittelbach Chris Rowley*

v1.2k 2019/11/14 printed October 26, 2020

This file is maintained by the L^AT_EX Project team.
Bug reports can be opened (category `latex`) at
<https://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	9
2.4	Loading Unicode mappings at begin document	14
3	Mapping characters — based on font (glyph) encodings	15
3.1	About the table itself	15
3.2	The mapping table	16
3.3	Notes	27
3.4	Mappings for OT1 glyphs	28
3.5	Mappings for OMS glyphs	28
3.6	Mappings for TS1 glyphs	28
3.7	Mappings for <code>latex.ltx</code> glyphs	29
3.8	Old <code>utf8.def</code> file as a temp fix for pT _E X and friends	29
4	A test document	34

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

1 Introduction

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

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.

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 <-utf8-2018> [2019/11/14 v1.2k UTF-8 support for inputenc]
16 <*utf8>

```

This is a temporary fix for the e-p \TeX / e-up \TeX engines that do not yet have a `\ifincname` primitive. Once this is available the extra file will be dropped.

```

17 \ifx\ifincname\undefined % old e-pTeX or e-upTeX engines
18 \input utf8-2018.def
19 \expandafter\@firstofone
20 \else
21 \expandafter\@gobble
22 \fi
23 \endinput
24 \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.

```

25 \catcode'\ \saved@space@catcode

```

2.2 Parsing UTF-8 input

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 remaining octets as the argument(s).

- When typesetting we pick up the necessary number of additional octets, check if they form a command that L^AT_EX knows about (via `\csname u8:\string #1\string #2...\endcsname`) and if so use that for typesetting. `\string` is needed as the octets may (all?) be active and we want the literal values in the name.
- If the UTF-8 character is going to be part of a label, then it is essentially becoming part of some csname and with the test `\ifincsname` we can find this out. If so, we render the whole sequence of octets harmless by using `\string` too when the starting octet executes (`\UTF@...@octets@string`).
- Another possible case is that `\protect` has *not* the meaning of `\typeset@protect`. In that case we may do a `\write` or we may do a `\protected@edef` or ... In all such cases we want to keep the sequence of octets unchanged, but we can't use `\string` this time, since at least in the case of `\protect@edef` the result may later be typeset after all (in fact that is quite likely) and so at that point the starting octet needs to be an active character again (the others could be stringified). So for this case we use `\noexpand` (`(\UTF@...s@octets@noexpand)`).

`\UTFviii@two@octets` Putting that all together the code for a start octet of a two byte sequence would then look like this:

```

26 \long\def\UTFviii@two@octets{%
27   \ifincsname
28     \expandafter \UTF@two@octets@string
29   \else
30     \ifx \protect\@typeset@protect \else
31       \expandafter\expandafter\expandafter \UTF@two@octets@noexpand
32     \fi
33   \fi
34 \UTFviii@two@octets@combine
35 }

```

`\ifincsname` is tested first because that can be true even if we are otherwise doing typesetting. If this is the case we use `\string` on the whole octet sequence. `\UTF@two@octets@string` not only does this but also gets rid of the command `\UTFviii@two@octets@combine` in the input stream by picking it up as a first argument and dropping it.

If this is not the case and we are doing typesetting (i.e., `\protect` is `\typeset@protect`), then we execute `\UTFviii@two@octets@combine` which picks up all octets and typesets the character (or generates an error if it doesn't know how to typeset it).

However, if we are not doing typesetting, then we execute the command `\UTFviii@two@octets@noexpand` which works like `\UTF@two@octets@string` but uses `\noexpand` instead of `\string`. This way the sequence is temporarily rendered harmless, e.g., would display as is or stays put inside a `\protected@edef`. But if the result is later reused the starting octet is still active and so will be able to construct the UTF-8 character again.

`\UTFviii@three@octets` The definitions for the other starting octets are the same except that they pick up
`\UTFviii@four@octets` more octets after them.

```

36 \long\def\UTFviii@three@octets{%

```

```

37 \ifincsname
38   \expandafter \UTF@three@octets@string
39 \else
40   \ifx \protect\@typeset@protect \else
41     \expandafter\expandafter\expandafter \UTF@three@octets@noexpand
42   \fi
43 \fi
44 \UTFviii@three@octets@combine
45 }

46 \long\def\UTFviii@four@octets{%
47   \ifincsname
48     \expandafter \UTF@four@octets@string
49   \else
50     \ifx \protect\@typeset@protect \else
51       \expandafter\expandafter\expandafter \UTF@four@octets@noexpand
52     \fi
53   \fi
54   \UTFviii@four@octets@combine
55 }

```

```

\UTFviii@two@octets@noexpand These temporarily prevent the active chars from expanding.
\UTFviii@three@octets@noexpand 56 \long\def\UTF@two@octets@noexpand#1#2{\noexpand#2\noexpand}
\UTFviii@four@octets@noexpand 57 \long\def\UTF@three@octets@noexpand#1#2#3{\noexpand#2\noexpand#3\noexpand}
58 \long\def\UTF@four@octets@noexpand#1#2#3#4{\noexpand#2\noexpand#3\noexpand#4\noexpand}

```

```

\UTFviii@two@octets@string And the same with \string for use in \csname constructions.
\UTFviii@three@octets@string 59 \long\def\UTF@two@octets@string#1#2{\string#2\string}
\UTFviii@four@octets@string 60 \long\def\UTF@three@octets@string#1#2#3{\string#2\string#3\string}
61 \long\def\UTF@four@octets@string#1#2#3#4{\string#2\string#3\string#4\string}

```

```

\UTFviii@two@octets@combine From the arguments a control sequence with a name of the form u8:#1#2... is
\UTFviii@three@octets@combine constructed where the #i (i > 1) are the arguments and #1 is the starting octet
\UTFviii@four@octets@combine (as a  $\TeX$  active character token). Since some or even all of these characters are
active 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 \LaTeX form).

```

62 \long\def\UTFviii@two@octets@combine#1#2{\expandafter
63   \UTFviii@defined\csname u8:\string#1\string#2\endcsname}

64 \long\def\UTFviii@three@octets@combine#1#2#3{\expandafter
65   \UTFviii@defined\csname u8:\string#1\string#2\string#3\endcsname}

66 \long\def\UTFviii@four@octets@combine#1#2#3#4{\expandafter
67   \UTFviii@defined\csname u8:\string#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.

```

```

68 \def\UTFviii@defined#1{%
69   \ifx#1\relax

```

Test if the sequence is invalid UTF-8 or valid UTF-8 but without a L^AT_EX definition.

```
70 \if\relax\expandafter\UTFviii@checkseq\string#1\relax\relax
```

The endline 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.

```
71 \UTFviii@undefined@err{#1}%
```

```
72 \else
```

```
73 \PackageError{inputenc}{Invalid UTF-8 byte sequence}%
```

```
74 \UTFviii@invalid@help
```

```
75 \fi
```

```
76 \else\expandafter
```

```
77 #1%
```

```
78 \fi
```

```
79 }
```

```
\UTFviii@invalid@err
```

```
\UTFviii@invalid@help
```

```
80 \def\UTFviii@invalid@err#1{%
```

```
81 \PackageError{inputenc}{Invalid UTF-8 byte "\UTFviii@hexnumber{'#1}}%
```

```
82 \UTFviii@invalid@help}
```

```
83 \def\UTFviii@invalid@help{%
```

```
84 The document does not appear to be in UTF-8 encoding.\MessageBreak
```

```
85 Try adding \noexpand\UseRawInputEncoding as the first line of the file.\MessageBreak
```

```
86 or specify an encoding such as \noexpand\usepackage[latin1]{inputenc}.\MessageBreak
```

```
87 in the document preamble.\MessageBreak
```

```
88 Alternatively, save the file in UTF-8 using your editor or another tool}
```

```
\UTFviii@undefined@err
```

```
89 \def\UTFviii@undefined@err#1{%
```

```
90 \PackageError{inputenc}{Unicode character \expandafter
```

```
91 \UTFviii@splitcsname\string#1\relax
```

```
92 \MessageBreak
```

```
93 not set up for use with LaTeX}%
```

```
94 {You may provide a definition with.\MessageBreak
```

```
95 \noexpand\DeclareUnicodeCharacter}%
```

```
96 }
```

```
\UTFviii@checkseq
```

Check that the csname consists of a valid UTF-8 sequence.

```
\UTFviii@check@continue
```

```
97 \def\UTFviii@checkseq#1:#2#3{%
```

```
98 \ifnum'#2<"80 %
```

```
99 \ifx\relax#3\else1\fi
```

```
100 \else
```

```
101 \ifnum'#2<"C0 %
```

```
102 1 %
```

```
103 \else
```

```
104 \expandafter\expandafter\expandafter\UTFviii@check@continue
```

```
105 \expandafter\expandafter\expandafter#3%
```

```
106 \fi
```

```
107 \fi}
```

```

108 \def\UTFviii@check@continue#1{%
109   \ifx\relax#1%
110   \else
111   \ifnum'#1<"80 1\else\ifnum'#1>"BF 1\fi\fi
112   \expandafter\UTFviii@check@continue
113   \fi
114 }

```

`\UTFviii@loop` This bit of code derived from `xmltex` defines the active character corresponding to starting octets to call `\UTFviii@two@octets` etc as appropriate. The starting octet itself is passed directly as the first argument, the others are picked up later en route.

The `\UTFviii@loop` loops through the numbers starting at `\count@` and ending 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.

```

115 \begingroup
116 \catcode'\~13
117 \catcode'\~12
118 \def\UTFviii@loop{%
119   \uccode'\~\count@
120   \uppercase\expandafter{\UTFviii@tmp}%
121   \advance\count@\@ne
122   \ifnum\count@<\@tempcnta
123   \expandafter\UTFviii@loop
124   \fi}

```

Handle the single byte control characters. C0 controls are valid UTF-8 but defined to give the “Character not defined error” They may be defined with `\DeclareUnicodeCharacter`.

```

125   \def\UTFviii@tmp{\xdef~{\noexpand\UTFviii@undefined@err{:\string~}}
126 % 0 ^^@ null
127   \count@"1
128   \@tempcnta9
129 % 9 ^^I tab
130 % 10 ^^J nl
131 \UTFviii@loop
132   \count@"11
133   \@tempcnta12
134 \UTFviii@loop
135 % 12 ^^L
136 % 13 ^^M
137   \count@"14
138   \@tempcnta32
139 \UTFviii@loop

```

Bytes with leading bits 10 are not valid UTF-8 starting bytes

```

140   \count@"80
141   \@tempcnta"C2
142   \def\UTFviii@tmp{\xdef~{\noexpand\UTFviii@invalid@err\string~}}
143 \UTFviii@loop

```

Setting up 2-byte UTF-8: The starting bytes is passed as an active character so that it can be reprocessed later!


```

144 \count@"C2
145 \@tempcnta"E0
146 \def\UTFviii@tmp{\xdef~{\noexpand\UTFviii@two@octets\noexpand~}}
147 \UTFviii@loop
    Setting up 3-byte UTF-8:
148 \count@"E0
149 \@tempcnta"F0
150 \def\UTFviii@tmp{\xdef~{\noexpand\UTFviii@three@octets\noexpand~}}
151 \UTFviii@loop
    Setting up 4-byte UTF-8:
152 \count@"F0
153 \@tempcnta"F5
154 \def\UTFviii@tmp{\xdef~{\noexpand\UTFviii@four@octets\noexpand~}}
155 \UTFviii@loop
    Bytes above F4 are not valid UTF-8 starting bytes as they would encode num-
bers beyond the Unicode range
156 \count@"F5
157 \@tempcnta"100
158 \def\UTFviii@tmp{\xdef~{\noexpand\UTFviii@invalid@err\string~}}
159 \UTFviii@loop
160 \endgroup

```

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

```
161 \@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.

```
162 \ifx\@begindocumenthook\@undefined
163 \makeatother
```

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

```
164 \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.

```

165 \begingroup
166 \catcode'\=12
167 \catcode'\<=12
168 \catcode'\.=12
169 \catcode'\,=12
170 \catcode'\;=12
171 \catcode'\!=12
172 \catcode'\~=13

```

```

173 \gdef\DeclareUnicodeCharacter#1#2{%
174   \count@"#1\relax
175   \wlog{ \space\space defining Unicode char U+#1 (decimal \the\count@)}%
176   \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).

```

177   \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.

```

178   \def\UTFviii@two@octets##1##2{\csname u8:##1\string##2\endcsname}%
179   \def\UTFviii@three@octets##1##2##3{\csname u8:##1%
180                                     \string##2\string##3\endcsname}%
181   \def\UTFviii@four@octets##1##2##3##4{\csname u8:##1%
182                                     \string##2\string##3\string##4\endcsname}%

```

Now we simply:-) need to use the right number of `\expandafter` 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

```

183   \expandafter\expandafter\expandafter
184   \expandafter\expandafter\expandafter
185   \expandafter
186   \gdef\UTFviii@tmp{\IeC{#2}}%
187   \endgroup
188 }

```

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

```

189 \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.

```

190 % \begingroup

```

The original code from `xmltex` 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).

```
191 % \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. In the case single byte UTF-8 sequences, only allow definition if the character is already active. The definition of `\UTFviii@tmp` looks slightly strange but is designed for the sequence of `\expandafter` in `\DeclareUnicodeCharacter`.

```
192 \ifnum\count@<"A0\relax
193   \ifnum\catcode\count@=13
194     \uccode'\~=\count@\uppercase{\def\UTFviii@tmp{\@empty\@empty}}%
195   \else
196     \PackageError{inputenc}%
197       {Cannot define non-active Unicode char value < 00A0}%
198     \eoha
199   \def\UTFviii@tmp{\UTFviii@tmp}%
200 \fi
```

The code below is derived from `xmltex` and generates the UTF-8 byte sequence for the number in `\count@`.

The reverse operation (just used in error messages) has now been added as `\decode@UTFviii`.

```
201 \else\ifnum\count@<"800\relax
202   \parse@UTFviii@a,%
203   \parse@UTFviii@b C\UTFviii@two@octets.,%
204 \else\ifnum\count@<"10000\relax
205   \parse@UTFviii@a;%
206   \parse@UTFviii@a,%
207   \parse@UTFviii@b E\UTFviii@three@octets.{,;%}
208 \else
```

Test added here for out of range values, the 4-octet definitions are still set up so that `\DeclareUnicodeCharacter` does something sensible if the user scrolls past this error.

```
209   \ifnum\count@>"10FFFF\relax
210     \PackageError{inputenc}%
211       {\UTFviii@hexnumber\count@\space too large for Unicode}%
212     {Values between 0 and 10FFFF are permitted}%
213   \fi
214   \parse@UTFviii@a;%
215   \parse@UTFviii@a,%
216   \parse@UTFviii@a!%
217   \parse@UTFviii@b F\UTFviii@four@octets.{!,;%}
218 \fi
219 \fi
220 \fi
221 % \endgroup
222 }
```

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

```
223 \gdef\parse@UTFviii@a#1{%
224     \@tempcnta\count@
225     \divide\count@ 64
226     \@tempcntb\count@
227     \multiply\count@ 64
228     \advance\@tempcnta-\count@
229     \advance\@tempcnta 128
230     \uccode'#1\@tempcnta
231     \count@\@tempcntb}
```

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

```
232 \gdef\parse@UTFviii@b#1#2#3#4{%
233     \advance\count@ "#10\relax
234     \uccode'#3\count@
235     \uppercase{\gdef\UTFviii@tmp{#2#3#4}}}
```

`\decode@UTFviii` In the reverse direction, take a sequence of octets(bytes) representing a character in UTF-8 and construct the Unicode number. The sequence is terminated by `\relax`.

In this version, if the sequence is not valid UTF-8 you probably get a low level arithmetic error from `\numexpr` or stray characters at the end. Getting a better error message would be somewhat expensive. As the main use is for reporting characters in messages, this is done just using expansion, so `\numexpr` is used, A stub returning 0 is defined if `\numexpr` is not available.

```
236 \ifx\numexpr\undefined
237 \gdef\decode@UTFviii#1{0}
238 \else
```

If the input is malformed UTF-8 there may not be enough closing `)` so add 5 so there are always some remaining then cleanup and remove any remaining ones at the end. This avoids `\numexpr` parse errors while outputting a package error.

```
239 \gdef\decode@UTFviii#1\relax{%
240     \expandafter\UTFviii@cleanup
241     \the\numexpr\dec@de@UTFviii#1\relax)))))\@empty}
242 \gdef\UTFviii@cleanup#1)#2\@empty{#1}
243 \gdef\dec@de@UTFviii#1{%
244 \ifx\relax#1%
245 \else
246 \ifnum'#1>"EF
247     (((('#1-"F0)%
248 \else
249 \ifnum'#1>"DF
250     (((('#1-"E0)%
251 \else
252 \ifnum'#1>"BF
253     ((''#1-"C0)%
254 \else
255 \ifnum'#1>"7F
256     )*64+('#1-"80)%
257 \else
```

```

258         +' #1 %
259         \fi
260     \fi
261     \fi
262 \fi
263 \expandafter\dec@de@UTFviii
264 \fi}
265 \fi

```

`\UTFviii@hexnumber` Convert a number to a sequence of uppercase hex digits. If `\numexpr` is not available, it returns its argument unchanged.

```

266 \ifx\numexpr\@undefined
267 \global\let\UTFviii@hexnumber\@firstofone
268 \global\UTFviii@hexdigit\hexnumber@
269 \else
270 \gdef\UTFviii@hexnumber#1{%
271 \ifnum#1>15 %
272 \expandafter\UTFviii@hexnumber\expandafter{\the\numexpr(#1-8)/16\relax}%
273 \fi
274 \UTFviii@hexdigit{\numexpr#1\ifnum#1>0-((#1-8)/16)*16\fi\relax}%
275 }

```

Almost but not quite `\hexnumber@`.

```

276 \gdef\UTFviii@hexdigit#1{\ifcase\numexpr#1\relax
277 0\or1\or2\or3\or4\or5\or6\or7\or8\or9\or
278 A\or B\or C\or D\or E\or F\fi}
279 \fi

```

`\UTFviii@splitcsname` Split a csname representing a unicode character and return the character and the
`\UTFviii@hexcodepoint` unicode number in hex.

```

280 \gdef\UTFviii@hexcodepoint#1{U+%
281 \ifnum#1<16 0\fi
282 \ifnum#1<256 0\fi
283 \ifnum#1<4096 0\fi
284 \UTFviii@hexnumber{#1}%
285 }%
286 \gdef\UTFviii@splitcsname#1:#2\relax{%

```

Need to pre-expand the argument to ensure cleanup in case of mal-formed UTF-8.

```

287 #2 (\expandafter\UTFviii@hexcodepoint\expandafter{%
288 \the\numexpr\decode@UTFviii#2\relax})%
289 }

```

```

290 \endgroup

```

```

291 \@onlypreamble\DeclareUnicodeCharacter

```

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

```

292 \@onlypreamble\parse@XML@charref
293 \@onlypreamble\parse@UTFviii@a
294 \@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 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.

```
295 \begingroup
296   \def\cdp@elt#1#2#3#4{%
297     \wlog{Now handling font encoding #1 ...}%
298     \lowercase{%
299       \InputIfFileExists{#1enc.dfu}}%
300       {\wlog{... processing UTF-8 mapping file for font %
301         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.

```
302       \catcode'\ 9\relax}%
303       {\wlog{... no UTF-8 mapping file for font encoding #1}}%
304   }
305   \cdp@list
306 \endgroup
```

However, we don't know if there are font encodings still to be loaded (either with `fontenc` or directly with `\input` by 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.

```
307 \def\DeclareFontEncoding@#1#2#3{%
308   \expandafter
309   \ifx\csname T@#1\endcsname\relax
310     \def\cdp@elt{\noexpand\cdp@elt}%
311     \xdef\cdp@list{\cdp@list\cdp@elt{#1}%
312       {\default@family}{\default@series}%
313       {\default@shape}}%
314     \expandafter\let\csname#1-cmd\endcsname\@changed@cmd
315     \begingroup
316       \wlog{Now handling font encoding #1 ...}%
317       \lowercase{%
318         \InputIfFileExists{#1enc.dfu}}%
319         {\wlog{... processing UTF-8 mapping file for font %
320           encoding #1}}%
321         {\wlog{... no UTF-8 mapping file for font encoding #1}}%
322       \endgroup
323   \else
324     \@font@info{Redefining font encoding #1}%
325   \fi
326   \global\@namedef{T@#1}{#2}%
327   \global\@namedef{M@#1}{\default@M#3}%
```

```

328 \xdef\LastDeclaredEncoding{#1}%
329 }
330 </utf8>

```

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 L^AT_EX 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 L^AT_EX 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 T_EX 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.

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, T_EX itself typically does not treat even many perfectly 'normal text' 7-bit slots as text characters, so it is unclear whether L^AT_EX 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 hex-digit number greater than 00BF and at most 10FFFF.

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

```
331 <all, t1, ot1, ly1>\DeclareUnicodeCharacter{00A0}{\nobreakspace}
332 <all, t1, ot1, ly1>\DeclareUnicodeCharacter{00A1}{\textexclamdown}
333 <all, ts1, ly1>\DeclareUnicodeCharacter{00A2}{\textcent}
334 <all, ts1, t1, ot1, ly1>\DeclareUnicodeCharacter{00A3}{\textsterling}
335 <all, x2, ts1, t2c, t2b, t2a, ly1, lcy>\DeclareUnicodeCharacter{00A4}{\textcurrency}
336 <all, ts1, ly1>\DeclareUnicodeCharacter{00A5}{\textyen}
337 <all, ts1, ly1>\DeclareUnicodeCharacter{00A6}{\textbrokenbar}
338 <all, x2, ts1, t2c, t2b, t2a, oms, ly1>\DeclareUnicodeCharacter{00A7}{\textsection}
339 <all, ts1>\DeclareUnicodeCharacter{00A8}{\textasciidieresis}
340 <all, ts1, utf8>\DeclareUnicodeCharacter{00A9}{\textcopyright}
341 <all, ts1, ly1, utf8>\DeclareUnicodeCharacter{00AA}{\textordfeminine}
342 <*all, x2, t2c, t2b, t2a, t1, ot2, ly1, lcy>
343 \DeclareUnicodeCharacter{00AB}{\guillemotleft}
344 </all, x2, t2c, t2b, t2a, t1, ot2, ly1, lcy>
345 <all, ts1>\DeclareUnicodeCharacter{00AC}{\textlnot}
346 <all, t1, ot1, ly1>\DeclareUnicodeCharacter{00AD}{\textdash}
347 <all, ts1, ly1, utf8>\DeclareUnicodeCharacter{00AE}{\textregistered}
348 <all, ts1>\DeclareUnicodeCharacter{00AF}{\textasciimacron}
349 <all, ts1, ly1>\DeclareUnicodeCharacter{00B0}{\textdegree}
350 <all, ts1>\DeclareUnicodeCharacter{00B1}{\textpm}
351 <all, ts1>\DeclareUnicodeCharacter{00B2}{\texttwosuperior}
352 <all, ts1>\DeclareUnicodeCharacter{00B3}{\textthreesuperior}
353 <all, ts1>\DeclareUnicodeCharacter{00B4}{\textasciiacute}
354 <all, ts1, ly1>\DeclareUnicodeCharacter{00B5}{\textmu} % micro sign
355 <all, ts1, oms, ly1>\DeclareUnicodeCharacter{00B6}{\textparagraph}
356 <all, oms, ts1, ly1>\DeclareUnicodeCharacter{00B7}{\textperiodcentered}
357 <all, ot1>\DeclareUnicodeCharacter{00B8}{\c }
358 <all, ts1>\DeclareUnicodeCharacter{00B9}{\textonesuperior}
359 <all, ts1, ly1, utf8>\DeclareUnicodeCharacter{00BA}{\textordmasculine}
360 <*all, x2, t2c, t2b, t2a, t1, ot2, ly1, lcy>
361 \DeclareUnicodeCharacter{00BB}{\guillemotright}
362 </all, x2, t2c, t2b, t2a, t1, ot2, ly1, lcy>
363 <all, ts1, ly1>\DeclareUnicodeCharacter{00BC}{\textonequarter}
364 <all, ts1, ly1>\DeclareUnicodeCharacter{00BD}{\textonehalf}
365 <all, ts1, ly1>\DeclareUnicodeCharacter{00BE}{\textthreequarters}
366 <all, t1, ot1, ly1>\DeclareUnicodeCharacter{00BF}{\textquestiondown}
367 <all, t1, ly1>\DeclareUnicodeCharacter{00C0}{\@tabacckludge'A}
368 <all, t1, ly1>\DeclareUnicodeCharacter{00C1}{\@tabacckludge'A}
369 <all, t1, ly1>\DeclareUnicodeCharacter{00C2}{\^A}
370 <all, t1, ly1>\DeclareUnicodeCharacter{00C3}{\~A}
371 <all, t1, ly1>\DeclareUnicodeCharacter{00C4}{\"A}
372 <all, t1, ot1, ly1>\DeclareUnicodeCharacter{00C5}{\r A}
373 <all, t1, ot1, ly1, lcy>\DeclareUnicodeCharacter{00C6}{\AE}
374 <all, t1, ly1>\DeclareUnicodeCharacter{00C7}{\c C}
375 <all, t1, ly1>\DeclareUnicodeCharacter{00C8}{\@tabacckludge'E}
376 <all, t1, ly1>\DeclareUnicodeCharacter{00C9}{\@tabacckludge'E}
377 <all, t1, ly1>\DeclareUnicodeCharacter{00CA}{\^E}
378 <all, t1, ly1>\DeclareUnicodeCharacter{00CB}{\"E}
379 <all, t1, ly1>\DeclareUnicodeCharacter{00CC}{\@tabacckludge'I}
```



```

380 <all, t1, ly1>\DeclareUnicodeCharacter{00CD}{\@tabacckludge'I}
381 <all, t1, ly1>\DeclareUnicodeCharacter{00CE}{\^I}
382 <all, t1, ly1>\DeclareUnicodeCharacter{00CF}{\~I}
383 <all, t1, ly1>\DeclareUnicodeCharacter{00D0}{\DH}
384 <all, t1, ly1>\DeclareUnicodeCharacter{00D1}{\~N}
385 <all, t1, ly1>\DeclareUnicodeCharacter{00D2}{\@tabacckludge'O}
386 <all, t1, ly1>\DeclareUnicodeCharacter{00D3}{\@tabacckludge'O}
387 <all, t1, ly1>\DeclareUnicodeCharacter{00D4}{\^O}
388 <all, t1, ly1>\DeclareUnicodeCharacter{00D5}{\~O}
389 <all, t1, ly1>\DeclareUnicodeCharacter{00D6}{\~O}
390 <all, ts1>\DeclareUnicodeCharacter{00D7}{\texttimes}
391 <all, t1, ot1, ly1, lcy>\DeclareUnicodeCharacter{00D8}{\O}
392 <all, t1, ly1>\DeclareUnicodeCharacter{00D9}{\@tabacckludge'U}
393 <all, t1, ly1>\DeclareUnicodeCharacter{00DA}{\@tabacckludge'U}
394 <all, t1, ly1>\DeclareUnicodeCharacter{00DB}{\~U}
395 <all, t1, ly1>\DeclareUnicodeCharacter{00DC}{\~U}
396 <all, t1, ly1>\DeclareUnicodeCharacter{00DD}{\@tabacckludge'Y}
397 <all, t1, ly1>\DeclareUnicodeCharacter{00DE}{\TH}
398 <all, t1, ot1, ly1, lcy>\DeclareUnicodeCharacter{00DF}{\ss}
399 <all, t1, ly1>\DeclareUnicodeCharacter{00E0}{\@tabacckludge'a}
400 <all, t1, ly1>\DeclareUnicodeCharacter{00E1}{\@tabacckludge'a}
401 <all, t1, ly1>\DeclareUnicodeCharacter{00E2}{\^a}
402 <all, t1, ly1>\DeclareUnicodeCharacter{00E3}{\~a}
403 <all, t1, ly1>\DeclareUnicodeCharacter{00E4}{\~a}
404 <all, t1, ly1>\DeclareUnicodeCharacter{00E5}{\r a}
405 <all, t1, ot1, ly1, lcy>\DeclareUnicodeCharacter{00E6}{\ae}
406 <all, t1, ly1>\DeclareUnicodeCharacter{00E7}{\c c}
407 <all, t1, ly1>\DeclareUnicodeCharacter{00E8}{\@tabacckludge'e}
408 <all, t1, ly1>\DeclareUnicodeCharacter{00E9}{\@tabacckludge'e}
409 <all, t1, ly1>\DeclareUnicodeCharacter{00EA}{\^e}
410 <all, t1, ly1>\DeclareUnicodeCharacter{00EB}{\~e}
411 <all, t1, ot1, ly1>\DeclareUnicodeCharacter{00EC}{\@tabacckludge'i}
412 <all, t1, ot1, ly1>\DeclareUnicodeCharacter{00ED}{\@tabacckludge'i}
413 <all, t1, ot1, ly1>\DeclareUnicodeCharacter{00EE}{\^i}
414 <all, t1, ot1, ly1>\DeclareUnicodeCharacter{00EF}{\~i}
415 <all, t1, ly1>\DeclareUnicodeCharacter{00F0}{\dh}
416 <all, t1, ly1>\DeclareUnicodeCharacter{00F1}{\~n}
417 <all, t1, ly1>\DeclareUnicodeCharacter{00F2}{\@tabacckludge'o}
418 <all, t1, ly1>\DeclareUnicodeCharacter{00F3}{\@tabacckludge'o}
419 <all, t1, ly1>\DeclareUnicodeCharacter{00F4}{\^o}
420 <all, t1, ly1>\DeclareUnicodeCharacter{00F5}{\~o}
421 <all, t1, ly1>\DeclareUnicodeCharacter{00F6}{\~o}
422 <all, ts1>\DeclareUnicodeCharacter{00F7}{\textdiv}
423 <all, t1, ot1, ly1, lcy>\DeclareUnicodeCharacter{00F8}{\o}
424 <all, t1, ly1>\DeclareUnicodeCharacter{00F9}{\@tabacckludge'u}
425 <all, t1, ly1>\DeclareUnicodeCharacter{00FA}{\@tabacckludge'u}
426 <all, t1, ly1>\DeclareUnicodeCharacter{00FB}{\~u}
427 <all, t1, ly1>\DeclareUnicodeCharacter{00FC}{\~u}
428 <all, t1, ly1>\DeclareUnicodeCharacter{00FD}{\@tabacckludge'y}
429 <all, t1, ly1>\DeclareUnicodeCharacter{00FE}{\th}
430 <all, t1, ly1>\DeclareUnicodeCharacter{00FF}{\~y}
431 <all, t1>\DeclareUnicodeCharacter{0100}{\@tabacckludge=A}
432 <all, t1>\DeclareUnicodeCharacter{0101}{\@tabacckludge=a}
433 <all, t1>\DeclareUnicodeCharacter{0102}{\u A}

```

434 <all,t1>\DeclareUnicodeCharacter{0103}{\u a}
435 <all,t1>\DeclareUnicodeCharacter{0104}{\k A}
436 <all,t1>\DeclareUnicodeCharacter{0105}{\k a}
437 <all,t1>\DeclareUnicodeCharacter{0106}{\@tabacckludge' C}
438 <all,t1>\DeclareUnicodeCharacter{0107}{\@tabacckludge' c}
439 <all,t1>\DeclareUnicodeCharacter{0108}{\^C}
440 <all,t1>\DeclareUnicodeCharacter{0109}{\^c}
441 <all,t1>\DeclareUnicodeCharacter{010A}{\ .C}
442 <all,t1>\DeclareUnicodeCharacter{010B}{\ .c}
443 <all,t1>\DeclareUnicodeCharacter{010C}{\v C}
444 <all,t1>\DeclareUnicodeCharacter{010D}{\v c}
445 <all,t1>\DeclareUnicodeCharacter{010E}{\v D}
446 <all,t1>\DeclareUnicodeCharacter{010F}{\v d}
447 <all,t1>\DeclareUnicodeCharacter{0110}{\DJ}
448 <all,t1>\DeclareUnicodeCharacter{0111}{\dj}
449 <all,t1>\DeclareUnicodeCharacter{0112}{\@tabacckludge=E}
450 <all,t1>\DeclareUnicodeCharacter{0113}{\@tabacckludge=e}
451 <all,t1>\DeclareUnicodeCharacter{0114}{\u E}
452 <all,t1>\DeclareUnicodeCharacter{0115}{\u e}
453 <all,t1>\DeclareUnicodeCharacter{0116}{\ .E}
454 <all,t1>\DeclareUnicodeCharacter{0117}{\ .e}
455 <all,t1>\DeclareUnicodeCharacter{0118}{\k E}
456 <all,t1>\DeclareUnicodeCharacter{0119}{\k e}
457 <all,t1>\DeclareUnicodeCharacter{011A}{\v E}
458 <all,t1>\DeclareUnicodeCharacter{011B}{\v e}
459 <all,t1>\DeclareUnicodeCharacter{011C}{\^G}
460 <all,t1>\DeclareUnicodeCharacter{011D}{\^g}
461 <all,t1>\DeclareUnicodeCharacter{011E}{\u G}
462 <all,t1>\DeclareUnicodeCharacter{011F}{\u g}
463 <all,t1>\DeclareUnicodeCharacter{0120}{\ .G}
464 <all,t1>\DeclareUnicodeCharacter{0121}{\ .g}
465 <all,t1>\DeclareUnicodeCharacter{0122}{\c G}
466 <all,t1>\DeclareUnicodeCharacter{0123}{\c g}
467 <all,t1>\DeclareUnicodeCharacter{0124}{\^H}
468 <all,t1>\DeclareUnicodeCharacter{0125}{\^h}
469 <all,t1>\DeclareUnicodeCharacter{0128}{\^I}
470 <all,t1>\DeclareUnicodeCharacter{0129}{\^i}
471 <all,t1>\DeclareUnicodeCharacter{012A}{\@tabacckludge=I}
472 <all,t1>\DeclareUnicodeCharacter{012B}{\@tabacckludge=i}
473 <all,t1>\DeclareUnicodeCharacter{012C}{\u I}
474 <all,t1>\DeclareUnicodeCharacter{012D}{\u i}
475 <all,t1>\DeclareUnicodeCharacter{012E}{\k I}

476 <all,t1>\DeclareUnicodeCharacter{012F}{\k i}
477 <all,t1>\DeclareUnicodeCharacter{0130}{\ .I}
478 <all,t2c,t2b,t2a,t1,ot2,ot1,ly1,lcy>\DeclareUnicodeCharacter{0131}{\i}
479 <all,t1>\DeclareUnicodeCharacter{0132}{\IJ}
480 <all,t1>\DeclareUnicodeCharacter{0133}{\ij}
481 <all,t1>\DeclareUnicodeCharacter{0134}{\^J}
482 <all,t1>\DeclareUnicodeCharacter{0135}{\^j}
483 <all,t1>\DeclareUnicodeCharacter{0136}{\c K}
484 <all,t1>\DeclareUnicodeCharacter{0137}{\c k}
485 <all,t1>\DeclareUnicodeCharacter{0139}{\@tabacckludge'L}
486 <all,t1>\DeclareUnicodeCharacter{013A}{\@tabacckludge'l}
487 <all,t1>\DeclareUnicodeCharacter{013B}{\c L}

488 <all, t1>\DeclareUnicodeCharacter{013C}{\c l}
489 <all, t1>\DeclareUnicodeCharacter{013D}{\v L}
490 <all, t1>\DeclareUnicodeCharacter{013E}{\v l}
491 <all, t1, ot1, ly1>\DeclareUnicodeCharacter{0141}{\L}
492 <all, t1, ot1, ly1>\DeclareUnicodeCharacter{0142}{\l}
493 <all, t1>\DeclareUnicodeCharacter{0143}{\@tabacckludge'N}
494 <all, t1>\DeclareUnicodeCharacter{0144}{\@tabacckludge'n}
495 <all, t1>\DeclareUnicodeCharacter{0145}{\c N}
496 <all, t1>\DeclareUnicodeCharacter{0146}{\c n}
497 <all, t1>\DeclareUnicodeCharacter{0147}{\v N}
498 <all, t1>\DeclareUnicodeCharacter{0148}{\v n}
499 <all, t1>\DeclareUnicodeCharacter{014A}{\NG}
500 <all, t1>\DeclareUnicodeCharacter{014B}{\ng}
501 <all, t1>\DeclareUnicodeCharacter{014C}{\@tabacckludge=O}
502 <all, t1>\DeclareUnicodeCharacter{014D}{\@tabacckludge=o}
503 <all, t1>\DeclareUnicodeCharacter{014E}{\u O}
504 <all, t1>\DeclareUnicodeCharacter{014F}{\u o}
505 <all, t1>\DeclareUnicodeCharacter{0150}{\H O}
506 <all, t1>\DeclareUnicodeCharacter{0151}{\H o}
507 <all, t1, ot1, ly1, lcy>\DeclareUnicodeCharacter{0152}{\OE}
508 <all, t1, ot1, ly1, lcy>\DeclareUnicodeCharacter{0153}{\oe}
509 <all, t1>\DeclareUnicodeCharacter{0154}{\@tabacckludge'R}
510 <all, t1>\DeclareUnicodeCharacter{0155}{\@tabacckludge'r}
511 <all, t1>\DeclareUnicodeCharacter{0156}{\c R}
512 <all, t1>\DeclareUnicodeCharacter{0157}{\c r}
513 <all, t1>\DeclareUnicodeCharacter{0158}{\v R}
514 <all, t1>\DeclareUnicodeCharacter{0159}{\v r}
515 <all, t1>\DeclareUnicodeCharacter{015A}{\@tabacckludge'S}
516 <all, t1>\DeclareUnicodeCharacter{015B}{\@tabacckludge's}
517 <all, t1>\DeclareUnicodeCharacter{015C}{\^S}
518 <all, t1>\DeclareUnicodeCharacter{015D}{\^s}
519 <all, t1>\DeclareUnicodeCharacter{015E}{\c S}
520 <all, t1>\DeclareUnicodeCharacter{015F}{\c s}
521 <all, t1, ly1>\DeclareUnicodeCharacter{0160}{\v S}
522 <all, t1, ly1>\DeclareUnicodeCharacter{0161}{\v s}
523 <all, t1>\DeclareUnicodeCharacter{0162}{\c T}
524 <all, t1>\DeclareUnicodeCharacter{0163}{\c t}
525 <all, t1>\DeclareUnicodeCharacter{0164}{\v T}
526 <all, t1>\DeclareUnicodeCharacter{0165}{\v t}
527 <all, t1>\DeclareUnicodeCharacter{0168}{\^U}
528 <all, t1>\DeclareUnicodeCharacter{0169}{\^u}
529 <all, t1>\DeclareUnicodeCharacter{016A}{\@tabacckludge=U}
530 <all, t1>\DeclareUnicodeCharacter{016B}{\@tabacckludge=u}
531 <all, t1>\DeclareUnicodeCharacter{016C}{\u U}
532 <all, t1>\DeclareUnicodeCharacter{016D}{\u u}
533 <all, t1>\DeclareUnicodeCharacter{016E}{\r U}
534 <all, t1>\DeclareUnicodeCharacter{016F}{\r u}
535 <all, t1>\DeclareUnicodeCharacter{0170}{\H U}
536 <all, t1>\DeclareUnicodeCharacter{0171}{\H u}
537 <all, t1>\DeclareUnicodeCharacter{0172}{\k U}
538 <all, t1>\DeclareUnicodeCharacter{0173}{\k u}

539 <all, t1, ot1, ly1>\DeclareUnicodeCharacter{0174}{\^W}
540 <all, t1, ot1, ly1>\DeclareUnicodeCharacter{0175}{\^w}
541 <all, t1, ot1, ly1>\DeclareUnicodeCharacter{0176}{\^Y}

```

542 <all, t1, ot1, ly1>\DeclareUnicodeCharacter{0177}{\~y}
543 <all, t1, ly1>\DeclareUnicodeCharacter{0178}{\~Y}
544 <all, t1>\DeclareUnicodeCharacter{0179}{\@tabacckludge'Z}
545 <all, t1>\DeclareUnicodeCharacter{017A}{\@tabacckludge'z}
546 <all, t1>\DeclareUnicodeCharacter{017B}{\ .Z}
547 <all, t1>\DeclareUnicodeCharacter{017C}{\ .z}
548 <all, t1, ly1>\DeclareUnicodeCharacter{017D}{\v Z}
549 <all, t1, ly1>\DeclareUnicodeCharacter{017E}{\v z}
550 <all, ts1, ly1>\DeclareUnicodeCharacter{0192}{\textflorin}

551 <all, t1>\DeclareUnicodeCharacter{01CD}{\v A}
552 <all, t1>\DeclareUnicodeCharacter{01CE}{\v a}
553 <all, t1>\DeclareUnicodeCharacter{01CF}{\v I}
554 <all, t1>\DeclareUnicodeCharacter{01D0}{\v \i}
555 <all, t1>\DeclareUnicodeCharacter{01D1}{\v O}
556 <all, t1>\DeclareUnicodeCharacter{01D2}{\v o}
557 <all, t1>\DeclareUnicodeCharacter{01D3}{\v U}
558 <all, t1>\DeclareUnicodeCharacter{01D4}{\v u}
559 <all, t1>\DeclareUnicodeCharacter{01E2}{\@tabacckludge=\AE}
560 <all, t1>\DeclareUnicodeCharacter{01E3}{\@tabacckludge=\ae}
561 <all, t1>\DeclareUnicodeCharacter{01E6}{\v G}
562 <all, t1>\DeclareUnicodeCharacter{01E7}{\v g}
563 <all, t1>\DeclareUnicodeCharacter{01E8}{\v K}
564 <all, t1>\DeclareUnicodeCharacter{01E9}{\v k}
565 <all, t1>\DeclareUnicodeCharacter{01EA}{\k O}
566 <all, t1>\DeclareUnicodeCharacter{01EB}{\k o}
567 <all, t1>\DeclareUnicodeCharacter{01F0}{\v \j}
568 <all, t1>\DeclareUnicodeCharacter{01F4}{\@tabacckludge'G}
569 <all, t1>\DeclareUnicodeCharacter{01F5}{\@tabacckludge'g}

570 <all, t1, ot1, ly1>\DeclareUnicodeCharacter{0218}{\textcommabelow S}
571 <all, t1, ot1, ly1>\DeclareUnicodeCharacter{0219}{\textcommabelow s}
572 <all, t1, ot1, ly1>\DeclareUnicodeCharacter{021A}{\textcommabelow T}
573 <all, t1, ot1, ly1>\DeclareUnicodeCharacter{021B}{\textcommabelow t}

574 <all, t1>\DeclareUnicodeCharacter{0232}{\@tabacckludge=Y}
575 <all, t1>\DeclareUnicodeCharacter{0233}{\@tabacckludge=y}
576 <all, t2c, t2b, t2a, t1, ot2, ot1, ly1, lcy>\DeclareUnicodeCharacter{0237}{\j}
577 <all, ly1, utf8>\DeclareUnicodeCharacter{02C6}{\textasciicircum}
578 <all, ts1>\DeclareUnicodeCharacter{02C7}{\textasciicaron}
579 <all, ly1, utf8>\DeclareUnicodeCharacter{02DC}{\textasciitilde}
580 <all, ts1>\DeclareUnicodeCharacter{02D8}{\textasciibreve}
581 <all, t1>\DeclareUnicodeCharacter{02D9}{\ .{}}
582 <all, t1>\DeclareUnicodeCharacter{02DB}{\k{}}
583 <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.

```

584 <*all, x2, t2c, t2b, t2a, ot2, lcy>
585 \DeclareUnicodeCharacter{0400}{\@tabacckludge'\CYRE}
586 </all, x2, t2c, t2b, t2a, ot2, lcy>
587 <all, x2, t2c, t2b, t2a, ot2, lcy>\DeclareUnicodeCharacter{0401}{\CYRYO}
588 <all, x2, t2a, ot2>\DeclareUnicodeCharacter{0402}{\CYRDJE}
589 <*all, x2, t2c, t2b, t2a, ot2, lcy>
590 \DeclareUnicodeCharacter{0403}{\@tabacckludge'\CYRG}
591 </all, x2, t2c, t2b, t2a, ot2, lcy>

```

592 <all, x2, t2a, ot2, lcy>\DeclareUnicodeCharacter{0404}{\CYRIE}
593 <all, x2, t2c, t2b, t2a, ot2>\DeclareUnicodeCharacter{0405}{\CYRDZE}
594 <all, x2, t2c, t2b, t2a, ot2, lcy>\DeclareUnicodeCharacter{0406}{\CYRII}
595 <all, x2, t2a, lcy>\DeclareUnicodeCharacter{0407}{\CYRYI}
596 <all, x2, t2c, t2b, t2a, ot2>\DeclareUnicodeCharacter{0408}{\CYRJE}
597 <all, x2, t2b, t2a, ot2>\DeclareUnicodeCharacter{0409}{\CYRLJE}
598 <all, x2, t2b, t2a, ot2>\DeclareUnicodeCharacter{040A}{\CYRNJE}
599 <all, x2, t2a, ot2>\DeclareUnicodeCharacter{040B}{\CYRTSHE}
600 <*all, x2, t2c, t2b, t2a, ot2, lcy>
601 \DeclareUnicodeCharacter{040C}{\@tabacckludge'\CYRK}
602 \DeclareUnicodeCharacter{040D}{\@tabacckludge'\CYRI}
603 </all, x2, t2c, t2b, t2a, ot2, lcy>
604 <all, x2, t2b, t2a, lcy>\DeclareUnicodeCharacter{040E}{\CYRUSHRT}
605 <all, x2, t2c, t2a, ot2>\DeclareUnicodeCharacter{040F}{\CYRDZHE}
606 <*all, x2, t2c, t2b, t2a, ot2, lcy>
607 \DeclareUnicodeCharacter{0410}{\CYRA}
608 \DeclareUnicodeCharacter{0411}{\CYRB}
609 \DeclareUnicodeCharacter{0412}{\CYRV}
610 \DeclareUnicodeCharacter{0413}{\CYRG}
611 \DeclareUnicodeCharacter{0414}{\CYRD}
612 \DeclareUnicodeCharacter{0415}{\CYRE}
613 \DeclareUnicodeCharacter{0416}{\CYRZH}
614 \DeclareUnicodeCharacter{0417}{\CYRZ}
615 \DeclareUnicodeCharacter{0418}{\CYRI}
616 \DeclareUnicodeCharacter{0419}{\CYRISHRT}
617 \DeclareUnicodeCharacter{041A}{\CYRK}
618 \DeclareUnicodeCharacter{041B}{\CYRL}
619 \DeclareUnicodeCharacter{041C}{\CYRM}
620 \DeclareUnicodeCharacter{041D}{\CYRN}
621 \DeclareUnicodeCharacter{041E}{\CYRO}
622 \DeclareUnicodeCharacter{041F}{\CYRP}
623 \DeclareUnicodeCharacter{0420}{\CYRR}
624 \DeclareUnicodeCharacter{0421}{\CYRS}
625 \DeclareUnicodeCharacter{0422}{\CYRT}
626 \DeclareUnicodeCharacter{0423}{\CYRU}
627 \DeclareUnicodeCharacter{0424}{\CYRF}
628 \DeclareUnicodeCharacter{0425}{\CYRH}
629 \DeclareUnicodeCharacter{0426}{\CYRC}
630 \DeclareUnicodeCharacter{0427}{\CYRCH}
631 \DeclareUnicodeCharacter{0428}{\CYRSH}
632 \DeclareUnicodeCharacter{0429}{\CYRSHCH}
633 \DeclareUnicodeCharacter{042A}{\CYRHRDSN}
634 \DeclareUnicodeCharacter{042B}{\CYRERY}
635 \DeclareUnicodeCharacter{042C}{\CYRSFTSN}
636 \DeclareUnicodeCharacter{042D}{\CYREREV}
637 \DeclareUnicodeCharacter{042E}{\CYRYU}
638 \DeclareUnicodeCharacter{042F}{\CYRYA}
639 \DeclareUnicodeCharacter{0430}{\cyr a}
640 \DeclareUnicodeCharacter{0431}{\cyr b}
641 \DeclareUnicodeCharacter{0432}{\cyr v}
642 \DeclareUnicodeCharacter{0433}{\cyr g}
643 \DeclareUnicodeCharacter{0434}{\cyr d}
644 \DeclareUnicodeCharacter{0435}{\cyr e}
645 \DeclareUnicodeCharacter{0436}{\cyr zh}

```

646 \DeclareUnicodeCharacter{0437}{\cyrz}
647 \DeclareUnicodeCharacter{0438}{\cyri}
648 \DeclareUnicodeCharacter{0439}{\cyrishrt}
649 \DeclareUnicodeCharacter{043A}{\cyrk}
650 \DeclareUnicodeCharacter{043B}{\cyr1}
651 \DeclareUnicodeCharacter{043C}{\cyrm}
652 \DeclareUnicodeCharacter{043D}{\cyrn}
653 \DeclareUnicodeCharacter{043E}{\cyro}
654 \DeclareUnicodeCharacter{043F}{\cyrp}
655 \DeclareUnicodeCharacter{0440}{\cyr}
656 \DeclareUnicodeCharacter{0441}{\cyr}
657 \DeclareUnicodeCharacter{0442}{\cyr}
658 \DeclareUnicodeCharacter{0443}{\cyr}
659 \DeclareUnicodeCharacter{0444}{\cyr}
660 \DeclareUnicodeCharacter{0445}{\cyr}
661 \DeclareUnicodeCharacter{0446}{\cyr}
662 \DeclareUnicodeCharacter{0447}{\cyr}
663 \DeclareUnicodeCharacter{0448}{\cyr}
664 \DeclareUnicodeCharacter{0449}{\cyr}
665 \DeclareUnicodeCharacter{044A}{\cyr}
666 \DeclareUnicodeCharacter{044B}{\cyr}
667 \DeclareUnicodeCharacter{044C}{\cyr}
668 \DeclareUnicodeCharacter{044D}{\cyr}
669 \DeclareUnicodeCharacter{044E}{\cyr}
670 \DeclareUnicodeCharacter{044F}{\cyr}
671 \DeclareUnicodeCharacter{0450}{\@tabacckludge'\cyre}
672 \DeclareUnicodeCharacter{0451}{\cyr}
673 </all, x2, t2c, t2b, t2a, ot2, lcy>
674 <all, x2, t2a, ot2>\DeclareUnicodeCharacter{0452}{\cyrdje}
675 <*all, x2, t2c, t2b, t2a, ot2, lcy>
676 \DeclareUnicodeCharacter{0453}{\@tabacckludge'\cyr}
677 </all, x2, t2c, t2b, t2a, ot2, lcy>
678 <all, x2, t2a, ot2, lcy>\DeclareUnicodeCharacter{0454}{\cyr}
679 <all, x2, t2c, t2b, t2a, ot2, lcy>\DeclareUnicodeCharacter{0455}{\cyrdze}
680 <all, x2, t2c, t2b, t2a, ot2, lcy>\DeclareUnicodeCharacter{0456}{\cyr}
681 <all, x2, t2a, lcy>\DeclareUnicodeCharacter{0457}{\cyr}
682 <all, x2, t2c, t2b, t2a, ot2>\DeclareUnicodeCharacter{0458}{\cyr}
683 <all, x2, t2b, t2a, ot2>\DeclareUnicodeCharacter{0459}{\cyr}
684 <all, x2, t2b, t2a, ot2>\DeclareUnicodeCharacter{045A}{\cyr}
685 <all, x2, t2a, ot2>\DeclareUnicodeCharacter{045B}{\cyr}
686 <*all, x2, t2c, t2b, t2a, ot2, lcy>
687 \DeclareUnicodeCharacter{045C}{\@tabacckludge'\cyr}
688 \DeclareUnicodeCharacter{045D}{\@tabacckludge'\cyr}
689 </all, x2, t2c, t2b, t2a, ot2, lcy>
690 <all, x2, t2b, t2a, lcy>\DeclareUnicodeCharacter{045E}{\cyr}
691 <all, x2, t2c, t2a, ot2>\DeclareUnicodeCharacter{045F}{\cyr}
692 <all, x2, ot2>\DeclareUnicodeCharacter{0462}{\CYRYAT}
693 <all, x2, ot2>\DeclareUnicodeCharacter{0463}{\cyr}
694 <all, x2>\DeclareUnicodeCharacter{046A}{\CYRBYUS}
695 <all, x2>\DeclareUnicodeCharacter{046B}{\cyr}

```

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).

```
696 <all, ot2>\DeclareUnicodeCharacter{0472}{\CYRFITA}
```

```
697 <all, ot2>\DeclareUnicodeCharacter{0473}{\cyrfita}
```

```
698 <all, x2, ot2>\DeclareUnicodeCharacter{0474}{\CYRIZH}
```

```
699 <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.

```
700 <all, x2>\DeclareUnicodeCharacter{0476}{\C\CYRIZH}
```

```
701 <all, x2>\DeclareUnicodeCharacter{0477}{\C\cyrizh}
```

```
702 <all, t2c>\DeclareUnicodeCharacter{048C}{\CYRSEMISFTSN}
```

```
703 <all, t2c>\DeclareUnicodeCharacter{048D}{\cyrsemisftsn}
```

```
704 <all, t2c>\DeclareUnicodeCharacter{048E}{\CYRRTICK}
```

```
705 <all, t2c>\DeclareUnicodeCharacter{048F}{\cyrrtick}
```

```
706 <all, x2, t2a, lcy>\DeclareUnicodeCharacter{0490}{\CYRGUP}
```

```
707 <all, x2, t2a, lcy>\DeclareUnicodeCharacter{0491}{\cyrgup}
```

```
708 <all, x2, t2b, t2a>\DeclareUnicodeCharacter{0492}{\CYRGHCRS}
```

```
709 <all, x2, t2b, t2a>\DeclareUnicodeCharacter{0493}{\cyrghcrs}
```

```
710 <all, x2, t2c, t2b>\DeclareUnicodeCharacter{0494}{\CYRGHK}
```

```
711 <all, x2, t2c, t2b>\DeclareUnicodeCharacter{0495}{\cyrghk}
```

```
712 <all, x2, t2b, t2a>\DeclareUnicodeCharacter{0496}{\CYRZHDSC}
```

```
713 <all, x2, t2b, t2a>\DeclareUnicodeCharacter{0497}{\cyrzhdsc}
```

```
714 <all, x2, t2a>\DeclareUnicodeCharacter{0498}{\CYRZDSC}
```

```
715 <all, x2, t2a>\DeclareUnicodeCharacter{0499}{\cyrzdsc}
```

```
716 <all, x2, t2c, t2b, t2a>\DeclareUnicodeCharacter{049A}{\CYRKDSC}
```

```
717 <all, x2, t2c, t2b, t2a>\DeclareUnicodeCharacter{049B}{\cyrkdsc}
```

```
718 <all, x2, t2a>\DeclareUnicodeCharacter{049C}{\CYRKVCRS}
```

```
719 <all, x2, t2a>\DeclareUnicodeCharacter{049D}{\cyrkvcrs}
```

```
720 <all, x2, t2c>\DeclareUnicodeCharacter{049E}{\CYRKHCRS}
```

```
721 <all, x2, t2c>\DeclareUnicodeCharacter{049F}{\cyrkhcrs}
```

```
722 <all, x2, t2a>\DeclareUnicodeCharacter{04A0}{\CYRKBEAK}
```

```
723 <all, x2, t2a>\DeclareUnicodeCharacter{04A1}{\cyrkbeak}
```

```
724 <all, x2, t2c, t2b, t2a>\DeclareUnicodeCharacter{04A2}{\CYRNDSC}
```

```
725 <all, x2, t2c, t2b, t2a>\DeclareUnicodeCharacter{04A3}{\cyrndsc}
```

```
726 <all, x2, t2b, t2a>\DeclareUnicodeCharacter{04A4}{\CYRNG}
```

```
727 <all, x2, t2b, t2a>\DeclareUnicodeCharacter{04A5}{\cyrng}
```

```
728 <all, x2, t2c>\DeclareUnicodeCharacter{04A6}{\CYRPHK}
```

```
729 <all, x2, t2c>\DeclareUnicodeCharacter{04A7}{\cyrphk}
```

```
730 <all, x2, t2c>\DeclareUnicodeCharacter{04A8}{\CYRABHHA}
```

```
731 <all, x2, t2c>\DeclareUnicodeCharacter{04A9}{\cyrabhha}
```

```
732 <all, x2, t2a>\DeclareUnicodeCharacter{04AA}{\CYRSDSC}
```

```
733 <all, x2, t2a>\DeclareUnicodeCharacter{04AB}{\cyrsdsc}
```

```
734 <all, x2, t2c>\DeclareUnicodeCharacter{04AC}{\CYRTDSC}
```

```
735 <all, x2, t2c>\DeclareUnicodeCharacter{04AD}{\cyrtdsc}
```

```
736 <all, x2, t2b, t2a>\DeclareUnicodeCharacter{04AE}{\CYRY}
```

```
737 <all, x2, t2b, t2a>\DeclareUnicodeCharacter{04AF}{\cyry}
```

```
738 <all, x2, t2a>\DeclareUnicodeCharacter{04B0}{\CYRYHCRS}
```

```

739 <all,x2,t2a>\DeclareUnicodeCharacter{04B1}{\cyryhcrs}
740 <all,x2,t2c,t2b,t2a>\DeclareUnicodeCharacter{04B2}{\CYRHDSC}
741 <all,x2,t2c,t2b,t2a>\DeclareUnicodeCharacter{04B3}{\cyrhds}
742 <all,x2,t2c>\DeclareUnicodeCharacter{04B4}{\CYRTETSE}
743 <all,x2,t2c>\DeclareUnicodeCharacter{04B5}{\cyrtetse}
744 <all,x2,t2c,t2b,t2a>\DeclareUnicodeCharacter{04B6}{\CYRCHRDSC}
745 <all,x2,t2c,t2b,t2a>\DeclareUnicodeCharacter{04B7}{\cyrchrdsc}
746 <all,x2,t2a>\DeclareUnicodeCharacter{04B8}{\CYRCHVCRS}
747 <all,x2,t2a>\DeclareUnicodeCharacter{04B9}{\cyrchvcrs}
748 <all,x2,t2c,t2b,t2a>\DeclareUnicodeCharacter{04BA}{\CYRSHHA}
749 <all,x2,t2c,t2b,t2a>\DeclareUnicodeCharacter{04BB}{\cyrshha}
750 <all,x2,t2c>\DeclareUnicodeCharacter{04BC}{\CYRABHCH}
751 <all,x2,t2c>\DeclareUnicodeCharacter{04BD}{\cyrabhch}
752 <all,x2,t2c>\DeclareUnicodeCharacter{04BE}{\CYRABHCHDSC}
753 <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.

```

754 <all,x2,t2c,t2b,t2a>\DeclareUnicodeCharacter{04C0}{\CYRpalochka}
755 <all,x2,t2c,t2b,t2a,ot2,lcy>\DeclareUnicodeCharacter{04C1}{\U\CYRZH}
756 <all,x2,t2c,t2b,t2a,ot2,lcy>\DeclareUnicodeCharacter{04C2}{\U\cyrz}
757 <all,x2,t2b>\DeclareUnicodeCharacter{04C3}{\CYRKHK}
758 <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?

```

759 <all,x2,t2c,t2b>\DeclareUnicodeCharacter{04C5}{\CYRLDSC}
760 <all,x2,t2c,t2b>\DeclareUnicodeCharacter{04C6}{\cyrl}
761 <all,x2,t2c,t2b>\DeclareUnicodeCharacter{04C7}{\CYRNHK}
762 <all,x2,t2c,t2b>\DeclareUnicodeCharacter{04C8}{\cyrn}
763 <all,x2,t2b>\DeclareUnicodeCharacter{04CB}{\CYRCHLDSC}
764 <all,x2,t2b>\DeclareUnicodeCharacter{04CC}{\cyrch}

```

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.

```

765 <all,x2,t2c>\DeclareUnicodeCharacter{04CD}{\CYRMDSC}
766 <all,x2,t2c>\DeclareUnicodeCharacter{04CE}{\cyr}
767 <all,x2,t2c,t2b,t2a,ot2,lcy>\DeclareUnicodeCharacter{04D0}{\U\CYRA}
768 <all,x2,t2c,t2b,t2a,ot2,lcy>\DeclareUnicodeCharacter{04D1}{\U\cyra}
769 <all,x2,t2c,t2b,t2a,ot2,lcy>\DeclareUnicodeCharacter{04D2}{\U\CYRA}
770 <all,x2,t2c,t2b,t2a,ot2,lcy>\DeclareUnicodeCharacter{04D3}{\U\cyra}
771 <all,x2,t2a>\DeclareUnicodeCharacter{04D4}{\CYRAE}
772 <all,x2,t2a>\DeclareUnicodeCharacter{04D5}{\cyrae}
773 <all,x2,t2c,t2b,t2a,ot2,lcy>\DeclareUnicodeCharacter{04D6}{\U\CYRE}
774 <all,x2,t2c,t2b,t2a,ot2,lcy>\DeclareUnicodeCharacter{04D7}{\U\cyre}
775 <all,x2,t2c,t2b,t2a>\DeclareUnicodeCharacter{04D8}{\CYRSCHWA}
776 <all,x2,t2c,t2b,t2a>\DeclareUnicodeCharacter{04D9}{\cyrschwa}
777 <all,x2,t2c,t2b,t2a>\DeclareUnicodeCharacter{04DA}{\U\CYRSCHWA}

```



```

778 <all,x2,t2c,t2b,t2a>\DeclareUnicodeCharacter{04DB}{\ "cyrschwa}
779 <all,x2,t2c,t2b,t2a,ot2,lcy>\DeclareUnicodeCharacter{04DC}{\ "CYRZH}
780 <all,x2,t2c,t2b,t2a,ot2,lcy>\DeclareUnicodeCharacter{04DD}{\ "cyrzh}
781 <all,x2,t2c,t2b,t2a,ot2,lcy>\DeclareUnicodeCharacter{04DE}{\ "CYRZ}
782 <all,x2,t2c,t2b,t2a,ot2,lcy>\DeclareUnicodeCharacter{04DF}{\ "cyrz}
783 <all,x2,t2c,t2b>\DeclareUnicodeCharacter{04E0}{\CYRABHDZE}
784 <all,x2,t2c,t2b>\DeclareUnicodeCharacter{04E1}{\cyrabhdze}
785 <all,x2,t2c,t2b,t2a,ot2,lcy>\DeclareUnicodeCharacter{04E2}{@tabacckludge=\CYRI}
786 <all,x2,t2c,t2b,t2a,ot2,lcy>\DeclareUnicodeCharacter{04E3}{@tabacckludge=\cyri}
787 <all,x2,t2c,t2b,t2a,ot2,lcy>\DeclareUnicodeCharacter{04E4}{\ "CYRI}
788 <all,x2,t2c,t2b,t2a,ot2,lcy>\DeclareUnicodeCharacter{04E5}{\ "cyri}
789 <all,x2,t2c,t2b,t2a,ot2,lcy>\DeclareUnicodeCharacter{04E6}{\ "CYRO}
790 <all,x2,t2c,t2b,t2a,ot2,lcy>\DeclareUnicodeCharacter{04E7}{\ "cyro}
791 <all,x2,t2c,t2b,t2a>\DeclareUnicodeCharacter{04E8}{\CYROTLD}
792 <all,x2,t2c,t2b,t2a>\DeclareUnicodeCharacter{04E9}{\cyrrotld}
793 <all,x2,t2c,t2b,t2a,ot2,lcy>\DeclareUnicodeCharacter{04EC}{\ "CYREREV}
794 <all,x2,t2c,t2b,t2a,ot2,lcy>\DeclareUnicodeCharacter{04ED}{\ "cyrerev}
795 <all,x2,t2c,t2b,t2a,ot2,lcy>\DeclareUnicodeCharacter{04EE}{@tabacckludge=\CYRU}
796 <all,x2,t2c,t2b,t2a,ot2,lcy>\DeclareUnicodeCharacter{04EF}{@tabacckludge=\cyru}
797 <all,x2,t2c,t2b,t2a,ot2,lcy>\DeclareUnicodeCharacter{04F0}{\ "CYRU}
798 <all,x2,t2c,t2b,t2a,ot2,lcy>\DeclareUnicodeCharacter{04F1}{\ "cyru}
799 <all,x2,t2c,t2b,t2a,ot2,lcy>\DeclareUnicodeCharacter{04F2}{\H\CYRU}
800 <all,x2,t2c,t2b,t2a,ot2,lcy>\DeclareUnicodeCharacter{04F3}{\H\cyru}
801 <all,x2,t2c,t2b,t2a,ot2,lcy>\DeclareUnicodeCharacter{04F4}{\ "CYRCH}
802 <all,x2,t2c,t2b,t2a,ot2,lcy>\DeclareUnicodeCharacter{04F5}{\ "cyrch}
803 <all,x2,t2b>\DeclareUnicodeCharacter{04F6}{\CYRGDSC}
804 <all,x2,t2b>\DeclareUnicodeCharacter{04F7}{\cyrgdsc}
805 <all,x2,t2c,t2b,t2a,ot2,lcy>\DeclareUnicodeCharacter{04F8}{\ "CYRERY}
806 <all,x2,t2c,t2b,t2a,ot2,lcy>\DeclareUnicodeCharacter{04F9}{\ "cyrery}
807 <all,t2b>\DeclareUnicodeCharacter{04FA}{\CYRGDSCHCRS}
808 <all,t2b>\DeclareUnicodeCharacter{04FB}{\cyrgdshcrs}
809 <all,x2,t2b>\DeclareUnicodeCharacter{04FC}{\CYRHHK}
810 <all,x2,t2b>\DeclareUnicodeCharacter{04FD}{\cyrhkh}
811 <all,t2b>\DeclareUnicodeCharacter{04FE}{\CYRHHCRS}
812 <all,t2b>\DeclareUnicodeCharacter{04FF}{\cyrhhcrs}
813 <all,ts1>\DeclareUnicodeCharacter{0E3F}{\textbaht}
814 <all,t1>\DeclareUnicodeCharacter{1E02}{\ .B}
815 <all,t1>\DeclareUnicodeCharacter{1E03}{\ .b}
816 <all,t1>\DeclareUnicodeCharacter{1E9E}{\SS}
817 <all,x2,t2c,t2b,t2a,t1,utf8>\DeclareUnicodeCharacter{200C}{\textcompwordmark}

818 <all,t1>\DeclareUnicodeCharacter{2010}{-}
819 <all,t1>\DeclareUnicodeCharacter{2011}{\mbox{-}}

```

U+2012 should be the width of a digit, endash is OK in many fonts including cm.

```

820 <all,t1>\DeclareUnicodeCharacter{2012}{\textendash}
821 <*all,x2,t2c,t2b,t2a,t1,ot2,ot1,ly1,lcy>
822 \DeclareUnicodeCharacter{2013}{\textendash}
823 \DeclareUnicodeCharacter{2014}{\textendash}

```

U+2015 is Horizontal bar

```

824 <all,t1>\DeclareUnicodeCharacter{2015}{\textendash}
825 </all,x2,t2c,t2b,t2a,t1,ot2,ot1,ly1,lcy>
826 <all,ts1>\DeclareUnicodeCharacter{2016}{\textbardbl}
827 <*all,x2,t2c,t2b,t2a,t1,ot2,ot1,lcy>

```

```

828 \DeclareUnicodeCharacter{2018}{\textquoteleft}
829 \DeclareUnicodeCharacter{2019}{\textquoteright}
830 \langle all, x2, t2c, t2b, t2a, t1, ot2, ot1, lcy)
831 \langle all, t1) \DeclareUnicodeCharacter{201A}{\quotesinglbase}
832 \langle *all, x2, t2c, t2b, t2a, t1, ot2, ot1, ly1, lcy)
833 \DeclareUnicodeCharacter{201C}{\textquotedblleft}
834 \DeclareUnicodeCharacter{201D}{\textquotedblright}
835 \langle all, x2, t2c, t2b, t2a, t1, ot2, ot1, ly1, lcy)
836 \langle all, x2, t2c, t2b, t2a, t1, lcy) \DeclareUnicodeCharacter{201E}{\quotedblbase}
837 \langle all, ts1, oms, ly1) \DeclareUnicodeCharacter{2020}{\textdagger}
838 \langle all, ts1, oms, ly1) \DeclareUnicodeCharacter{2021}{\textdaggerdbl}
839 \langle all, ts1, oms, ly1) \DeclareUnicodeCharacter{2022}{\textbullet}
840 \langle all, ly1, utf8) \DeclareUnicodeCharacter{2026}{\textellipsis}
841 \langle *all, x2, ts1, t2c, t2b, t2a, t1, ly1)
842 \DeclareUnicodeCharacter{2030}{\textperthousand}
843 \langle all, x2, ts1, t2c, t2b, t2a, t1, ly1)
844 \langle *all, x2, ts1, t2c, t2b, t2a, t1)
845 \DeclareUnicodeCharacter{2031}{\textpertenthousand}
846 \langle all, x2, ts1, t2c, t2b, t2a, t1)
847 \langle all, t1, ly1) \DeclareUnicodeCharacter{2039}{\guilsinglleft}
848 \langle all, t1, ly1) \DeclareUnicodeCharacter{203A}{\guilsinglright}
849 \langle all, ts1) \DeclareUnicodeCharacter{203B}{\textreferencemark}
850 \langle all, ts1) \DeclareUnicodeCharacter{203D}{\textinterrobang}
851 \langle all, ts1) \DeclareUnicodeCharacter{2044}{\textfractionsolidus}
852 \langle all, ts1) \DeclareUnicodeCharacter{204E}{\textasteriskcentered}
853 \langle all, ts1) \DeclareUnicodeCharacter{2052}{\textdiscount}
854 \langle all, ts1) \DeclareUnicodeCharacter{20A1}{\textcolonmonetary}
855 \langle all, ts1) \DeclareUnicodeCharacter{20A4}{\textlira}
856 \langle all, ts1) \DeclareUnicodeCharacter{20A6}{\textnaira}
857 \langle all, ts1) \DeclareUnicodeCharacter{20A9}{\textwon}
858 \langle all, ts1) \DeclareUnicodeCharacter{20AB}{\textdong}
859 \langle all, ts1) \DeclareUnicodeCharacter{20AC}{\texteuro}
860 \langle all, ts1) \DeclareUnicodeCharacter{20B1}{\textpeso}
861 \langle all, ts1) \DeclareUnicodeCharacter{2103}{\textcelsius}
862 \langle all, x2, ts1, t2c, t2b, t2a, ot2, lcy) \DeclareUnicodeCharacter{2116}{\textnumero}
863 \langle all, ts1) \DeclareUnicodeCharacter{2117}{\textcircledP}
864 \langle all, ts1) \DeclareUnicodeCharacter{211E}{\textrecipe}
865 \langle all, ts1) \DeclareUnicodeCharacter{2120}{\textservicemark}
866 \langle all, ts1, ly1, utf8) \DeclareUnicodeCharacter{2122}{\texttrademark}
867 \langle all, ts1) \DeclareUnicodeCharacter{2126}{\textohm}
868 \langle all, ts1) \DeclareUnicodeCharacter{2127}{\textmho}
869 \langle all, ts1) \DeclareUnicodeCharacter{212E}{\textestimated}
870 \langle all, ts1) \DeclareUnicodeCharacter{2190}{\textleftarrow}
871 \langle all, ts1) \DeclareUnicodeCharacter{2191}{\textuparrow}
872 \langle all, ts1) \DeclareUnicodeCharacter{2192}{\textrightarrow}
873 \langle all, ts1) \DeclareUnicodeCharacter{2193}{\textdownarrow}
874 \langle all, x2, ts1, t2c, t2b, t2a) \DeclareUnicodeCharacter{2329}{\textlangle}
875 \langle all, x2, ts1, t2c, t2b, t2a) \DeclareUnicodeCharacter{232A}{\textrangle}
876 \langle all, ts1) \DeclareUnicodeCharacter{2422}{\textblank}
877 \langle all, x2, t2c, t2b, t2a, t1, utf8) \DeclareUnicodeCharacter{2423}{\textvisiblespace}
878 \langle all, ts1) \DeclareUnicodeCharacter{25E6}{\textopenbullet}
879 \langle all, ts1) \DeclareUnicodeCharacter{25EF}{\textbigcircle}
880 \langle all, ts1) \DeclareUnicodeCharacter{266A}{\textmusicalnote}
881 \langle all, x2, ts1, t2c, t2b, t2a) \DeclareUnicodeCharacter{27E8}{\textlangle}

```

```

882 <all,x2,ts1,t2c,t2b,t2a>\DeclareUnicodeCharacter{27E9}{\texttriangle}
883 <all,t1>\DeclareUnicodeCharacter{1E20}{\@tabacckludge=G}
884 <all,t1>\DeclareUnicodeCharacter{1E21}{\@tabacckludge=g}

```

When doing cut-and-paste from other documents f-ligatures might show up as Unicode characters. We translate them back to individual characters so that they get accepted. If supported by the font (which is normally the case) they are then reconstructed as ligatures so they come out as desired. Otherwise they will come out as individual characters which is fine too.

```

885 <all,t1,ot1,ly1,t2a,t2b,t2c>\DeclareUnicodeCharacter{FB00}{ff} % ff
886 <all,t1,ot1,ly1,t2a,t2b,t2c>\DeclareUnicodeCharacter{FB01}{fi} % fi
887 <all,t1,ot1,ly1,t2a,t2b,t2c>\DeclareUnicodeCharacter{FB02}{fl} % fl
888 <all,t1,ot1,ly1,t2a,t2b,t2c>\DeclareUnicodeCharacter{FB03}{ffi} % ffi
889 <all,t1,ot1,ly1,t2a,t2b,t2c>\DeclareUnicodeCharacter{FB04}{ffl} % ffl
890 <all,t1,ot1,ly1,t2a,t2b,t2c>\DeclareUnicodeCharacter{FB05}{st} % st -- this is the long s (not
891 <all,t1,ot1,ly1,t2a,t2b,t2c>\DeclareUnicodeCharacter{FB06}{st} % st
892 <all,ts1,utf8>\DeclareUnicodeCharacter{FEFF}{\ifhmode\nobreak\fi}

```

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 T_EX 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.

```
893 <*utf8>
894 % This space is intentionally empty ...
895 </utf8>
```

3.8 Old utf8.def file as a temp fix for pTeX and friends

```
896 <*utf8-2018>
897 \ProvidesFile{utf8.def}
898 [2018/10/05 v1.2f UTF-8 support for inputenc]
899 \makeatletter
900 \catcode'\ \saved@space@catcode
901 \long\def\UTFviii@two@octets#1#2{\expandafter
902   \UTFviii@defined\csname u8:#1\string#2\endcsname}
903 \long\def\UTFviii@three@octets#1#2#3{\expandafter
904   \UTFviii@defined\csname u8:#1\string#2\string#3\endcsname}
905 \long\def\UTFviii@four@octets#1#2#3#4{\expandafter
906   \UTFviii@defined\csname u8:#1\string#2\string#3\string#4\endcsname}
907 \def\UTFviii@defined#1{%
908   \ifx#1\relax
909     \if\relax\expandafter\UTFviii@checkseq\string#1\relax\relax
910     \UTFviii@undefined@err{#1}%
911     \else
912     \PackageError{inputenc}{Invalid UTF-8 byte sequence}%
913       \UTFviii@invalid@help
914     \fi
915   \else\expandafter
916     #1%
917   \fi
918 }
919 \def\UTFviii@invalid@err#1{%
920   \PackageError{inputenc}{Invalid UTF-8 byte "\UTFviii@hexnumber{‘#1}}%
921     \UTFviii@invalid@help}
922 \def\UTFviii@invalid@help{%
923   The document does not appear to be in UTF-8 encoding.\MessageBreak
924   Try adding \noexpand\UseRawInputEncoding as the first line of the file\MessageBreak
925   or specify an encoding such as \noexpand\usepackage[latin1]{inputenc}\MessageBreak
926   in the document preamble.\MessageBreak
927   Alternatively, save the file in UTF-8 using your editor or another tool}
928 \def\UTFviii@undefined@err#1{%
929   \PackageError{inputenc}{Unicode character \expandafter
930     \UTFviii@splitcsname\string#1\relax
931     \MessageBreak
932     not set up for use with LaTeX}%
933     {You may provide a definition with\MessageBreak
934     \noexpand\DeclareUnicodeCharacter}%
935   }
936 \def\UTFviii@checkseq#1:#2#3{%
```

```

937 \ifnum'#2<"80 %
938 \ifx\relax#3\else1\fi
939 \else
940 \ifnum'#2<"C0 %
941 1 %
942 \else
943 \expandafter\expandafter\expandafter\UTFviii@check@continue
944 \expandafter\expandafter\expandafter#3%
945 \fi
946 \fi}
947 \def\UTFviii@check@continue#1{%
948 \ifx\relax#1%
949 \else
950 \ifnum'#1<"80 1\else\ifnum'#1>"BF 1\fi\fi
951 \expandafter\UTFviii@check@continue
952 \fi
953 }
954 \begingroup
955 \catcode'\~13
956 \catcode'\~12
957 \def\UTFviii@loop{%
958 \uccode'\~\count@
959 \uppercase\expandafter{\UTFviii@tmp}%
960 \advance\count@\@ne
961 \ifnum\count@<@\tempcnta
962 \expandafter\UTFviii@loop
963 \fi}
964 \def\UTFviii@tmp{\xdef~{\noexpand\UTFviii@undefined@err{:~string~}}
965 \count@"1
966 \@tempcnta9
967 \UTFviii@loop
968 \count@"11
969 \@tempcnta12
970 \UTFviii@loop
971 \count@"14
972 \@tempcnta32
973 \UTFviii@loop
974 \count@"80
975 \@tempcnta"C2
976 \def\UTFviii@tmp{\xdef~{\noexpand\UTFviii@invalid@err~string~}}
977 \UTFviii@loop
978 \count@"C2
979 \@tempcnta"E0
980 \def\UTFviii@tmp{\xdef~{\noexpand\UTFviii@two@octets~string~}}
981 \UTFviii@loop
982 \count@"E0
983 \@tempcnta"F0
984 \def\UTFviii@tmp{\xdef~{\noexpand\UTFviii@three@octets~string~}}
985 \UTFviii@loop
986 \count@"F0
987 \@tempcnta"F5
988 \def\UTFviii@tmp{\xdef~{\noexpand\UTFviii@four@octets~string~}}
989 \UTFviii@loop
990 \count@"F5

```

```

991 \@tempcnta"100
992 \def\UTFviii@tmp{\xdef~{\noexpand\UTFviii@invalid@err\string~}}
993 \UTFviii@loop
994 \endgroup
995 \@inpenc@test
996 \ifx\@begindocumenthook\@undefined
997 \makeatother
998 \endinput \fi
999 \begingroup
1000 \catcode'\ "=12
1001 \catcode'\ <=12
1002 \catcode'\ .=12
1003 \catcode'\ ,=12
1004 \catcode'\ ;=12
1005 \catcode'\ !=12
1006 \catcode'\ ^=13
1007 \gdef\DeclareUnicodeCharacter#1#2{%
1008 \count@"#1\relax
1009 \wlog{ \space\space defining Unicode char U+#1 (decimal \the\count@)}%
1010 \begingroup
1011 \parse@XML@charref
1012 \def\UTFviii@two@octets##1##2{\csname u8:##1\string##2\endcsname}%
1013 \def\UTFviii@three@octets##1##2##3{\csname u8:##1%
1014 \string##2\string##3\endcsname}%
1015 \def\UTFviii@four@octets##1##2##3##4{\csname u8:##1%
1016 \string##2\string##3\string##4\endcsname}%
1017 \expandafter\expandafter\expandafter
1018 \expandafter\expandafter\expandafter
1019 \expandafter
1020 \gdef\UTFviii@tmp{\IeC{#2}}%
1021 \endgroup
1022 }
1023 \gdef\parse@XML@charref{%
1024 \ifnum\count@<"A0\relax
1025 \ifnum\catcode\count@=13
1026 \uccode'\ ^=\count@\uppercase{\def\UTFviii@tmp{\@empty\@empty}}}%
1027 \else
1028 \PackageError{inputenc}%
1029 {Cannot define non-active Unicode char value < 00A0}%
1030 \@eha
1031 \def\UTFviii@tmp{\UTFviii@tmp}%
1032 \fi
1033 \else\ifnum\count@<"800\relax
1034 \parse@UTFviii@a,%
1035 \parse@UTFviii@b C\UTFviii@two@octets.,%
1036 \else\ifnum\count@<"10000\relax
1037 \parse@UTFviii@a;%
1038 \parse@UTFviii@a,%
1039 \parse@UTFviii@b E\UTFviii@three@octets.{,;}%
1040 \else
1041 \ifnum\count@>"10FFFF\relax
1042 \PackageError{inputenc}%
1043 {\UTFviii@hexnumber\count@\space too large for Unicode}%
1044 {Values between 0 and 10FFFF are permitted}%

```

```

1045     \fi
1046     \parse@UTFviii@a;%
1047     \parse@UTFviii@a,%
1048     \parse@UTFviii@a!%
1049     \parse@UTFviii@b F\UTFviii@four@octets.{!,,}%
1050     \fi
1051     \fi
1052     \fi
1053 }
1054 \gdef\parse@UTFviii@a#1{%
1055     \@tempcnta\count@
1056     \divide\count@ 64
1057     \@tempcntb\count@
1058     \multiply\count@ 64
1059     \advance\@tempcnta-\count@
1060     \advance\@tempcnta 128
1061     \uccode'#1\@tempcnta
1062     \count@\@tempcntb}
1063 \gdef\parse@UTFviii@b#1#2#3#4{%
1064     \advance\count@ "#10\relax
1065     \uccode'#3\count@
1066     \uppercase{\gdef\UTFviii@tmp{#2#3#4}}
1067 \ifx\numexpr\undefined
1068 \gdef\decode@UTFviii#1{0}
1069 \else
1070 \gdef\decode@UTFviii#1\relax{%
1071     \expandafter\UTFviii@cleanup
1072     \the\numexpr\dec@de@UTFviii#1\relax)))))\@empty}
1073 \gdef\UTFviii@cleanup#1)#2\@empty{#1}
1074 \gdef\dec@de@UTFviii#1{%
1075 \ifx\relax#1%
1076 \else
1077     \ifnum'#1>"EF
1078         (((('#1-"F0)%
1079     \else
1080         \ifnum'#1>"DF
1081             (((('#1-"E0)%
1082         \else
1083             \ifnum'#1>"BF
1084                 ((''#1-"C0)%
1085             \else
1086                 \ifnum'#1>"7F
1087                     )*64+(''#1-"80)%
1088                 \else
1089                     +'#1 %
1090                 \fi
1091             \fi
1092         \fi
1093     \fi
1094     \expandafter\dec@de@UTFviii
1095 \fi}
1096 \fi
1097 \ifx\numexpr\undefined
1098 \global\let\UTFviii@hexnumber\@firstofone

```



```

1099 \global\UTFviii@hexdigit\hexnumber@
1100 \else
1101 \gdef\UTFviii@hexnumber#1{%
1102 \ifnum#1>15 %
1103 \expandafter\UTFviii@hexnumber\expandafter{\the\numexpr(#1-8)/16\relax}%
1104 \fi
1105 \UTFviii@hexdigit{\numexpr#1\ifnum#1>0-((#1-8)/16)*16\fi\relax}%
1106 }
1107 \gdef\UTFviii@hexdigit#1{\ifcase\numexpr#1\relax
1108 0\or1\or2\or3\or4\or5\or6\or7\or8\or9\or
1109 A\or B\or C\or D\or E\or F\fi}
1110 \fi
1111 \gdef\UTFviii@hexcodepoint#1{U+%
1112 \ifnum#1<16 0\fi
1113 \ifnum#1<256 0\fi
1114 \ifnum#1<4096 0\fi
1115 \UTFviii@hexnumber{#1}%
1116 }%
1117 \gdef\UTFviii@splitcsname#1:#2\relax{%
1118 #2 (\expandafter\UTFviii@hexcodepoint\expandafter{%
1119 \the\numexpr\decode@UTFviii#2\relax})%
1120 }
1121 \endgroup
1122 \@onlypreamble\DeclareUnicodeCharacter
1123 \@onlypreamble\parse@XML@charref
1124 \@onlypreamble\parse@UTFviii@a
1125 \@onlypreamble\parse@UTFviii@b
1126 \begingroup
1127 \def\cdp@elt#1#2#3#4{%
1128 \wlog{Now handling font encoding #1 ...}%
1129 \lowercase{%
1130 \InputIfFileExists{#1enc.dfu}}%
1131 {\wlog{... processing UTF-8 mapping file for font %
1132 encoding #1}%
1133 \catcode'\ 9\relax}%
1134 {\wlog{... no UTF-8 mapping file for font encoding #1}}%
1135 }
1136 \cdp@list
1137 \endgroup
1138 \def\DeclareFontEncoding@#1#2#3{%
1139 \expandafter
1140 \ifx\csname T@#1\endcsname\relax
1141 \def\cdp@elt{\noexpand\cdp@elt}%
1142 \xdef\cdp@list{\cdp@list\cdp@elt{#1}%
1143 {\default@family}{\default@series}%
1144 {\default@shape}}%
1145 \expandafter\let\csname#1-cmd\endcsname\@changed@cmd
1146 \begingroup
1147 \wlog{Now handling font encoding #1 ...}%
1148 \lowercase{%
1149 \InputIfFileExists{#1enc.dfu}}%
1150 {\wlog{... processing UTF-8 mapping file for font %
1151 encoding #1}}%
1152 {\wlog{... no UTF-8 mapping file for font encoding #1}}%

```

```

1153     \endgroup
1154 \else
1155     \font@info{Redeclaring font encoding #1}%
1156 \fi
1157 \global\@namedef{T@#1}{#2}%
1158 \global\@namedef{M@#1}{\default@M#3}%
1159 \xdef\LastDeclaredEncoding{#1}%
1160 }
1161 \DeclareUnicodeCharacter{00A9}{\textcopyright}
1162 \DeclareUnicodeCharacter{00AA}{\textordfeminine}
1163 \DeclareUnicodeCharacter{00AE}{\textregistered}
1164 \DeclareUnicodeCharacter{00BA}{\textordmasculine}
1165 \DeclareUnicodeCharacter{02C6}{\textasciicircum}
1166 \DeclareUnicodeCharacter{02DC}{\textasciitilde}
1167 \DeclareUnicodeCharacter{200C}{\textcompwordmark}
1168 \DeclareUnicodeCharacter{2026}{\textellipsis}
1169 \DeclareUnicodeCharacter{2122}{\texttrademark}
1170 \DeclareUnicodeCharacter{2423}{\textvisiblespace}
1171 \DeclareUnicodeCharacter{FEFF}{\ifhmode\nobreak\fi}
1172 \endinput
1173 </utf8-2018>

```

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.

```

1174 <*test>
1175 \documentclass{article}
1176
1177 \usepackage[latin1,utf8]{inputenc}
1178 \usepackage[T1]{fontenc}
1179 \usepackage{trace}
1180
1181 \scrollmode % to run past the error below
1182
1183 \begin{document}
1184
1185 German umlauts in UTF-8: ^^c3^^a4^^c3^^b6^^c3^^bc %% äü
1186
1187 \inputencoding{latin1} % switch to latin1
1188
1189 German umlauts in UTF-8 but read by latin1 (and will produce one
1190 error since \verb=\textcurrency= is not provided):
1191 ^^c3^^a4^^c3^^b6^^c3^^bc
1192
1193 \inputencoding{utf8} % switch back to utf8
1194
1195 German umlauts in UTF-8: ^^c3^^a4^^c3^^b6^^c3^^bc
1196
1197
1198 Some codes that should produce errors as nothing is set up
1199 for them: ^^c3F ^^e1^^a4^^b6
1200
1201 And some that are not legal utf8 sequences: ^^c3X ^^e1XY

```

```
1202
1203 \showoutput
1204 \tracingstats=2
1205 \stop
1206 </test>
```