

The `luainputenc` package

Manuel Pégourié-Gonnard `mpg@elzevir.fr`
Élie Roux `elie.roux@telecom-bretagne.eu`

2010/10/05 v0.971

Abstract

Input encoding management for LuaTEX, needed only for compatibility with old documents. For new documents, using UTF-8 encoding and Unicode fonts is *strongly* recommended. You've been warned!

Contents

1 Overview: When (not) to use this package	1
2 Documentation	2
2.1 Introduction	2
2.2 Overview of 8-bit mode	2
2.3 Overview of UTF-8 mode	3
2.3.1 legacy mode	3
2.3.2 unicode font mode	4
2.3.3 mixed mode	4
3 Accessing the encoding in lua	4
4 Files	4
4.1 <code>inputenc.sty</code> patch	4
4.2 <code>luainputenc.sty</code>	5
4.3 <code>lutf8.def</code>	11
4.4 <code>lutf8x.def</code>	13
4.5 <code>luainputenc.lua</code>	16
5 Test file	22

1 Overview: When (not) to use this package

This package is strictly meant for compatibility. It is usefull in the two (overlapping) following cases:

1. Your source is not encoded in UTF-8 and you don't want to reencode it for some reason.

2. Your document is using legacy 8-bit fonts (with `fontenc`), as opposed to modern Unicode fonts (most probably with `fontspec` or `luatoflode` and `fontenc` with option EU2).

Surprisingly enough, in the second case `luainputenc` is needed, due to the way L^AT_EX implements font encodings.

From the user point of view, adapting an old document for LuaT_EX is really easy: replacing `inputenc` by `luainputenc` in the preamble is enough.

Note that `luainputenc` automatically loads `inputenc` if called with an old engine, so you will still be able to compile your documents with pdfT_EX without changing them.

`luainputenc` has several modes of operation. By default, it basically turns LuaT_EX into an 8-bit engine, which means you loose half of the benefits from using LuaT_EX. If you are using only Unicode fonts, you can activate a nicer mode of operation using the `unactivate` package option. That way, LuaT_EX remains a true Unicode engine.

Unicode fonts with LuaT_EX are handled using a new encoding: EU2. It is used internally by the `fontspec` package when loading Unicode fonts. This encoding is special as it needs non-ASCII characters to be non-active (unlike other font encodings), so you cannot mix old encodings and EU2. If you're using only Unicode fonts, this isn't a problem: use the `unactivate` package option mentioned in the previous paragraph.

But if you want to use both 8-bit fonts and Unicode fonts in your document, you need to use another package option, `lutf8x`. This option overrides L^AT_EX's mechanism for font encoding switching, so that it (un)activates non-ASCII characters on-the-fly. With this options, you'll be able change the font encoding from/to EU2, for example:

```
abc
{
\fontencoding{EU2}\usefont
\font\foo="MyOtfFont.otf"\foo
abc
}
abc
```

2 Documentation

2.1 Introduction

One of the most interesting new features of LuaT_EX is the fact that it is (like Omega/Aleph) not limited to 256 characters, and can now understand Unicode. The problem is that it does not read input the way older engines (like pdfT_EX) do, and thus `inputenc` is totally broken with LuaT_EX. This package aims at replacing `inputenc` for LuaT_EX, by adapting the way LuaT_EX handles input, and the way `inputenc` handles UTF-8. This package has two very distinct modes: 8-bit and UTF-8.

2.2 Overview of 8-bit mode

This package **does not** map 8-bit encodings to utf8. It allows LuaT_EX to read 8-bit characters, by converting each byte into a unicode character with the same character number. The

resulting unicode characters are not true UTF-8, they are what we will call “fake UTF-8”. For example the byte 225 will be converted into the unicode character with number 225 (two bytes long). It will be true UTF-8 only if the encoding is latin1.

Here is how it works: the 8-bit encodings are converted into fake UTF-8, so that the corresponding tokens are chars with the good numbers. Then (like `inputenc`) it reads the char numbers, and converts it into LICR (LATEX Internal Character Representation), with the font encoding.

In LATEX version 0.43, a new callback called `process_output_buffer`, this callbacks allows to make LATEX write 8-bit instead of UTF-8, so the behaviour is the same as pdfTeX as this level. For versions prior to 0.43 though, we need to do more tricky things, described in the next paragraph. This machinery is disabled for LATEX version 0.43 and superior, so you can keep the default behaviour, which will be compatible with pdfTeX in most cases, but you can consider the machinery obsolete.

For these old versions, `luainputenc` only changes the input behaviour, it does not change the ouput behaviour (when files are written for example). The consequence is that files will still be written by LATEX in UTF-8 (fake UTF-8 in this case), even if the asked input encoding is a 8-bit encoding. In most cases it’s not a problem, as most files will be written in LICR, meaning ASCII, which is both 8-bit and UTF-8. The problem comes when characters with a number > 128 are written in a 8-bit encoding. This may happen if you use `\protect` in a section for example. In these cases, LATEX will write fake UTF-8, and try to read 8-bit encoding, so it will get confused.

The proposed solution is to unactivate the input conversion when we read certain files or extention. This package should work with no change for most documents, but if you cook your own aux files with an unknown extention, you may have to force the package to read some files in UTF-8 instead of 8-bit. See comments in the `.sty` file to know the useful commands.

2.3 Overview of UTF-8 mode

The behaviour of `inputenc` in `utf8` mode is to read the input byte by byte, and decide if the character we are in is 1, 2, 3 or 4 bytes long, and then read other bytes accordingly. This behaviour fails with LATEX because it reads input character by character (characters do not have a fixed number of bytes in unicode). The result is thus an error.

All characters recognized by TEX are active characters, that correspond to a LICR macro. Then `inputenc` reads the `*.dfu` files that contain the correspondance between these LICR macros and a character number in the fonts for different font encodings (T1, OT1, etc.).

2.3.1 legacy mode

`luainputenc` can get this behaviour (we will call it *legacy mode*, but another difference implied by the fact that LATEX can read more than 256 characters is that fonts can also have more than 256 characters. LATEX can thus read unicode fonts. If we want to use unicode fonts (OTF for example), we can’t use the *legacy mode* anymore, as it would mean that we would have to rewrite a specially long `unicode.dfu` file, and it would be totally inefficient, as for instance é (unicode character number 233) would be mapped to `\'e`, and then mapped back to `\char 233`.

2.3.2 unicode font mode

To fix this, the most simple solution is to deactivate all activated characters, thus typing é will directly call \char 233 in the unicode fonts, and produce a é. We will call this behaviour the *unicode font mode*. To enable this mode, you can use the option `unactivate` in `luainputenc`, and you must use the font encoding EU2 provided by the `euenc` package. See documentation of `euenc` package for more details about EU2. To use this mode with EU2, you must be able to open OTF fonts. A simple way to do so it by using the package `luaotfload`.

2.3.3 mixed mode

But the *unicode font mode* has a strong limitation (that will certainly dissapear with time): it cannot use non-unicode fonts. If you want to mix unicode fonts and old fonts, you'll have to use the *mixed mode*. In this mode you can type some parts of your document in *legacy mode* and some in *unicode font mode*. The reason why we chose not to integrate this choice in the *legacy mode* is that we wanted to have a mode that preserved most of the backward compatibility, to safely compile old documents; the *mixed mode* introduces new things that may break old documents. To get the *mixed mode*, you must pass the option `lutf8x` to `luainputenc`. This mode is the most experimental.

3 Accessing the encoding in lua

In order to access the encoding and the package option in lua, two variables are set: `luainputenc.package_option` contains the option passed to the package, and `luainputenc.encoding` that contains the encoding (defaults to utf8, and is utf8 even with the options `unactivate`, `utf8x`, etc.).

4 Files

This package contains a .sty file for both L^AT_EX and Plain, a patch for inputenc to use `luainputenc` so that you can process old documents without changing anything, and the lua functions.

4.1 `inputenc.sty` patch

A good thing would be to patch `inputenc` to load `luainputenc` instead, so that you don't have to change your documents to load `luainputenc` especially. The L^AT_EX team is extremely conservative and does not want this patch applied (maybe we will find a solution later). Here is a patch for `inputenc.sty`:

```
1
2   \ifnum\@tempcnta<'#2\relax
3     \advance\@tempcnta\@ne
4     \repeat}
5 +
6 +\begingroup\expandafter\expandafter\expandafter\endgroup
7 +\expandafter\ifx\csname XeTeXversion\endcsname\relax\else
```

```

8 + \RequirePackage{xetex-inputenc}
9 + \DeclareOption*{\PassOptionsToPackage{\CurrentOption}{xetex-inputenc}}
10 + \ProcessOptions*
11 + \expandafter\endinput
12 +\fi
13 +\begingroup\expandafter\expandafter\expandafter\endgroup
14 +\expandafter\ifx\csname directlua\endcsname\relax\else
15 + \RequirePackage{luainputenc}
16 + \DeclareOption*{\PassOptionsToPackage{\CurrentOption}{luainputenc}}
17 + \ProcessOptions*
18 + \expandafter\endinput
19 +\fi
20 +
21 \ProcessOptions
22 \endinput
23 %%
24

```

4.2 luainputenc.sty

This file has some code from `inputenc.sty`, but also provides new options, and new macros to convert from 8-bit to fake UTF-8.

```

25 %
26 %% This file was adapted from inputenc.sty, which copyright is:
27 %% Copyright 1993 1994 1995 1996 1997 1998 1999 2000 2001 2002 2003 2004
28 %% 2005 2006 The LaTeX3 Project.
29 %%
30 %% inputenc.sty is under the lppl version 1.3c or later, and can be
31 %% found in the base LaTeX system.
32 %%
33 %% The lppl can be found at http://www.latex-project.org/lppl.txt
34 %%
35 %% The changes to inputenc.sty are Copyright 2009 Elie Roux, and are
36 %% under the CCO license.
37 %%
38 %% The changes are LuaTeX support.
39 %%
40 %% This file is distributed under the CCO license, with clause 6 of the
41 %% lppl as additional restrictions.
42

```

First we check if we are called with `LuaTeX`, `(pdf)TeX` or `XeTeX`. If we are called with `pdfTeX`, we default to `inputenc`, and to `xetex-inputenc` if we are called with `XeTeX`. We also remap the new options to `utf8` in these cases.

```

43
44 \RequirePackage{ifluatex}
45 \RequirePackage{ifxetex}
46
47 \ifxetex
48   \DeclareOption{unactivate}{\PassOptionsToPackage{utf8}{xetex-inputenc}}

```

```

49 \DeclareOption{lutf8}{\PassOptionsToPackage{utf8}{xetex-inputenc}}
50 \DeclareOption{lutf8x}{\PassOptionsToPackage{utf8}{xetex-inputenc}}
51 \DeclareOption*{\PassOptionsToPackage{\CurrentOption}{xetex-inputenc}}
52 \ProcessOptions*
53 \RequirePackage{xetex-inputenc}
54 \expandafter\endinput
55 \fi
56
57 \ifluatex\else
58 \DeclareOption{unactivate}{\PassOptionsToPackage{utf8}{inputenc}}
59 \DeclareOption{lutf8}{\PassOptionsToPackage{utf8}{inputenc}}
60 \DeclareOption{lutf8x}{\PassOptionsToPackage{utf8}{inputenc}}
61 \DeclareOption*{\PassOptionsToPackage{\CurrentOption}{inputenc}}
62 \ProcessOptions*
63 \RequirePackage{inputenc}
64 \expandafter\endinput
65 \fi
66

```

Here we know we are called with `LuATEX`. We first require `luatextra`, then we load the `lua` file.

```

67
68 \RequirePackage{luatexbase}
69 \RequirePackage{luatexbase-mcb}
70 \RequirePackage{luatexbase-modutils}
71
72 \luatexUseModule{luainputenc}
73

```

Here is some code from `inputenc`.

```

74
75 \def\DeclareInputMath#1{%
76   \inpc@t{#1}%
77   \bgroup
78     \uccode`~#1%
79     \uppercase{%
80       \egroup
81       \def~{%
82         }%
83     }%
84 \def\DeclareInputText#1#2{%
85   \def\reserved@a##1 ${}%
86   \def\reserved@b{#2}%
87   \ifcat\_ \expandafter\reserved@a\meaning\reserved@b$ ${}_%
88     \def\reserved@a{\meaning\reserved@b$ ${}_%
89   \else
90     \def\reserved@a{\IeC{#2}}%
91   \fi
92 }%
93 \def\IeC{%

```

```

94 \ifx\protect\@typeset@protect
95   \expandafter\@firstofone
96 \else
97   \noexpand\IeC
98 \fi
99 }

```

We changed a little the behaviour of this macro: we removed `\@inpenc@loop\^\^\?\^\^ff`, because it made no sense in UTF-8 mode. We will call this line for 8-bit encodings.

Note that the code has been changed for `\endlinechar`, because in new versions (from v0.43) of LuaTeX the value cannot exceed 127. Thus, with the old version of `luainputenc`, when trying to add 10000, it fails silently, and when 10000 is subtracted, the new value is -1, resulting in no end of lines at all in the document.

```

100
101 \def\inputencoding#1{%
102   \the\inpenc@prehook
103   \gdef\@inpenc@test{\global\let\@inpenc@test\relax}%
104   \edef\@inpenc@undefined{\noexpand\@inpenc@undefined@{#1}}%
105   \edef\inputencodingname{#1}%
106   \@inpenc@loop\^\^\A\^\^H%
107   \@inpenc@loop\^\^\K\^\^K%
108   \@inpenc@loop\^\^\N\^\^_%
109   \xdef\@saved@endlinechar{\the\endlinechar }%
110   \endlinechar=-1
111   \xdef\@saved@space@\catcode{\the\catcode`\ }%
112   \catcode`\ 9\relax
113   \input{#1.def}%
114   \endlinechar=\@saved@endlinechar{}%
115   \catcode`\ \@saved@space@\catcode\relax
116   \ifx\@inpenc@test\relax\else
117     \PackageWarning{inputenc}%
118       {No characters defined\MessageBreak
119        by input encoding change to '#1'\MessageBreak}%
120   \fi
121   \the\inpenc@posthook
122   \luadirect[luainputenc.set_option([#1])]%
123 }
124 \newtoks\inpenc@prehook
125 \newtoks\inpenc@posthook
126 \def\@inpenc@undefined@#1{\PackageError{inputenc}%
127   {Keyboard character used is undefined\MessageBreak
128    in inputencoding '#1'}%
129   {You need to provide a definition with
130    \noexpand\DeclareInputText\MessageBreak or
131    \noexpand\DeclareInputMath before using this key.}%
132 \def\@inpenc@loop#1#2{%
133   \tempcnta`#1\relax
134   \loop
135     \catcode\tempcnta\active
136     \bgroup

```

```

137      \uccode`~\@tempcnta
138      \uppercase{%
139      \egroup
140      \let~\inpc@undefined
141      }%
142 \ifnum\@tempcnta<#2\relax
143   \advance\@tempcnta\@ne
144 \repeat}
145

```

Here we declare our options. Note that we remap `utf8` to `lutf8`, because we use out `lutf8.def` instead of `inputenc`'s `utf8.def`.

```

146
147 \DeclareOption{utf8}{%
148   \inputencoding{lutf8}%
149 }
150
151 \DeclareOption{lutf8}{%
152   \inputencoding{lutf8}%
153 }
154
155 \DeclareOption{utf8x}{%
156   \inputencoding{lutf8}%
157 }
158
159 \DeclareOption{lutf8x}{%
160   \inputencoding{lutf8x}%
161 }
162

```

For the `unactivate` option, for *unicode font mode*, we just don't do anything.

```

163
164 \DeclareOption{unactivate}{%
165   \edef\inputencodingname{unactivate}%
166   \luadirect{luainputenc.set_option([[unactivate]])}
167 }
168

```

All other options are 8-bit encodings, so we activate the translation into fake UTF-8, and we execute the loop we removes from `\inputencoding`.

```

169
170 \DeclareOption*{%
171   \lIE@activate %
172   \inpc@loop\^\^\?\\^ff%
173   \inputencoding{\CurrentOption}%
174 }
175

```

The rest of the file is only the machinery for LuaTeX versions without the callback `process_output_buffer`, so it will be deprecated after TeXLive 2009, you are not advised to use it.

```

176
177 \ifnum\luatexversion>42
178
179   \newcommand*\lIE@activate}[0]{%
180     \luadirect{luainputenc.register_callbacks()}%
181   }
182
183 \else
184

\lIE@setstarted and \lIE@setstopped are called when the fake UTF-8 translation
must be activated or deactivated. You can call them several successive times. They are
called very often, even if the package is not activated (for example if it's loaded with the
utf8 option), but they act only if the package is activated.

185
186 \newcommand*\lIE@setstarted[0]{%
187   \ifnum\lIE@activated=1 %
188     \luadirect{luainputenc.setstarted()}%
189   \fi %
190 }
191
192 \newcommand*\lIE@setstopped[0]{%
193   \ifnum\lIE@activated=1 %
194     \luadirect{luainputenc.setstopped()}%
195   \fi %
196 }
197

```

The following 5 macros are made to declare a file that will have to be read in fake UTF-8 and not in 8-bit. These files are the ones that will be generated by TeX. In **no way** this means you can include true UTF-8 files, it means that you can include files that have been written by LuaTeX with luainputenc, which means files in fake UTF-8. The macros are very simple, when you call them with a file name (the same as the one you will use with \input), it will read it with or without the fake UTF-8 translation. This package includes a whole bunch of extention that will be read in fake UTF-8, so the occasions to use these macros will be rare, but if you use them, please report it to the package maintainer.

\lIE@SetUtfFile If you call this macro with a file name, each time you will input this file, it will be read in fake UTF-8. You can call it with a file that you generate with LuaTeX and that you want to include.

```

198
199 \newcommand*\lIE@SetUtfFile[1]{%
200   \luadirect{luainputenc.set_unicode_file([[#1]])}%
201 }
202

```

\lIE@SetNonUtfFile Same as the previous macro, except that the file will be read as 8-bit. This macro is useful if there is an exception in an extention (see further comments).

203

```
204 \newcommand*\LIE@SetNonUtfFile[1]{%
205   \luadirect{luainputenc.set_non_unicode_file([[#1]])}%
206 }
207
```

\LIE@UnsetFile This macro gives a file the default behaviour of its extention.

```
208
209 \newcommand*\LIE@UnsetFile[1]{%
210   \luadirect{luainputenc.unset_file([[#1]])}%
211 }
212
```

\LIE@SetUtfExt You can tell luainputenc to treat all files with a particular extention in a certain way. The way the file extention is checked is to compare the four last characters of the filename. So if your extention has only three letters, you must include the preceding dot. This macro tells luainputenc to read all files from an extention in fake UTF-8.

```
213
214 \newcommand*\LIE@SetUtfExt[1]{%
215   \luadirect{luainputenc.set_unicode_extention([[#1]])}%
216 }
217
```

\LIE@SetUtfExt Same as before, but the files will be read in 8-bit.

```
218
219 \newcommand*\LIE@SetNonUtfExt[1]{%
220   \luadirect{luainputenc.set_non_unicode_extention([[#1]])}%
221 }
222
```

\LIE@InputUtfFile This macro inputs a file in fake UTF-8. It has the "feature" to unset the behaviour on the file you will call, so to be safe, you must call them with files for which the behaviour has not been set.

```
223
224
225 \newcommand*\LIE@InputUtfFile[1]{%
226   \LIE@SetUtfFile{#1}%
227   \input #1%
228   \LIE@UnsetFile{#1}%
229 }
230
```

\LIE@InputNonUtfFile Same as before, but to read a file as 8-bit.

```
231
232 \newcommand*\LIE@InputNonUtfFile[1]{%
233   \LIE@SetNonUtfFile{#1}%
234   \input #1%
235   \LIE@UnsetFile{#1}%
236 }
237
```

Two definitions to put the previous two macros in the user space.

```
238 \newcommand*\InputUtfFile[1]{%
239   \lIE@InputUtfFile{#1}%
240 }
241
242
243 \newcommand*\InputNonUtfFile[1]{%
244   \lIE@InputNonUtfFile{#1}%
245 }
246
247 \newcount\lIE@activated
248
249 \newcommand*{\lIE@activate}[0]{%
250   \lIE@activated=1 %
251   \lIE@setstarted %
252 }
253
254 \newcommand*{\lIE@FromInputenc}[1]{%
255   \ifnum\lIE@activated=0 %
256     \lIE@activate %
257   \fi%
258 }
259
260 \fi
261
262 \ProcessOptions*
263
```

4.3 lutf8.def

```
264 %% This file was adapted from utf8.def, which copyright is:
265 %% Copyright 1993 1994 1995 1996 1997 1998 1999 2000 2001 2002 2003
266 %% 2004 2005 2006 The LaTeX3 Project.
267 %%
268 %% utf8.def is under the lppl version 1.3c or later, and can be found
269 %% in the base LaTeX system.
270 %%
271 %% The lppl can be found at http://www.latex-project.org/lppl.txt
272 %%
273 %% The changes to utf8.def are Copyright 2009 Elie Roux, and are under
274 %% the CCO license.
275 %%
276 %% The changes are LuaTeX support.
277 %%
278 %% This file is distributed under the CCO license, with clause 6 of the
279 %% lppl as additional restrictions.
280
```

Most of the file is taken from `utf8.def`, the main changes are commented. A lot of code was removed, especially the codes that analysed the unicode characters byte by byte.

```

281
282
283 \ProvidesFile{utf8.def}
284     [2010/05/10 v0.97 UTF-8 support for luainputenc]
285
286 \makeatletter
287 \catcode`\ \saved@space@catcode
288
289 \inplace@test
290
291 \ifx\@begindocumenthook\undefined
292     \makeatother
293     \endinput \fi
294

```

This function is changed a lot. Its aim is to map the character (first argument) to a macro (second argument). In `utf8.def` it was complicated as unicode was analyzed byte by byte. With `LuaTeX` it is extremely simple, we just have to activate the character, and call a traditional `\DeclareInputText`.

```

295
296 \gdef\DeclareUnicodeCharacter#1#2{%
297     \tempcnta"##1%
298     \catcode\tempcnta\active %
299     \DeclareInputText{\the\tempcnta}{##2}%
300 }
301
302 \onlypreamble\DeclareUnicodeCharacter
303
304 \def\cdp@elt#1#2#3#4{%
305     \wlog{Now handling font encoding #1 ...}%
306     \lowercase{%
307         \IfFileExists{#1enc.dfu}%
308             {\wlog{... processing UTF-8 mapping file for font encoding
309                 ##1}%
310             \catcode`\ 9\relax}%
311             {\wlog{... no UTF-8 mapping file for font encoding ##1}}%
312 }
313 \cdp@list
314
315 \def\DeclareFontEncoding@#1#2#3{%
316     \expandafter %
317     \ifx\csname T@#1\endcsname\relax %
318         \def\cdp@elt{\noexpand\cdp@elt}%
319         \xdef\cdp@list{\cdp@list\cdp@elt{#1}%
320             {\default@family}{\default@series}%
321             {\default@shape}}%
322         \expandafter\let\csname#1-cmd\endcsname\@changed@cmd %
323         \begin{group} %
324             \wlog{Now handling font encoding #1 ...}%
325             \lowercase{%
326                 \IfFileExists{#1enc.dfu}}%

```

```

327          {\wlog{... processing UTF-8 mapping file for font encoding #1}}%
328          {\wlog{... no UTF-8 mapping file for font encoding #1}}%
329      \endgroup
330  \else
331      \font@info{Redeclaring font encoding #1}%
332  \fi
333  \global\@namedef{T@#1}{#2}%
334  \global\@namedef{M@#1}{\default@M#3}%
335  \xdef\LastDeclaredEncoding{#1}%
336 }
337
338 \DeclareUnicodeCharacter{00A9}{\textcopyright}
339 \DeclareUnicodeCharacter{00AA}{\textordfeminine}
340 \DeclareUnicodeCharacter{00AE}{\textregistered}
341 \DeclareUnicodeCharacter{00BA}{\textordmASCULINE}
342 \DeclareUnicodeCharacter{02C6}{\textasciicircum}
343 \DeclareUnicodeCharacter{02DC}{\textasciitilde}
344 \DeclareUnicodeCharacter{200C}{\textcompwordmark}
345 \DeclareUnicodeCharacter{2026}{\textellipsis}
346 \DeclareUnicodeCharacter{2122}{\texttrademark}
347 \DeclareUnicodeCharacter{2423}{\textvisiblespace}
348

```

4.4 lutf8x.def

```

349 %% This file was adapted from utf8.def, which copyright is:
350 %% Copyright 1993 1994 1995 1996 1997 1998 1999 2000 2001 2002 2003
351 %% 2004 2005 2006 The LaTeX3 Project.
352 %%
353 %% utf8.def is under the lppl version 1.3c or later, and can be found
354 %% in the base LaTeX system.
355 %%
356 %% The lppl can be found at http://www.latex-project.org/lppl.txt
357 %%
358 %% The changes to utf8.def are Copyright 2009 Elie Roux, and are under
359 %% the CCO license.
360 %%
361 %% The changes are LuaTeX support.
362 %%
363 %% This file is distributed under the CCO license, with clause 6 of the
364 %% lppl as additional restrictions.
365

```

This file is mostly the code from `lutf.def`, but it adds mechanisms to pass from *legacy mode* to *unicode font mode*. The trick is to put in a lua table all characters that are activated by the *legacy mode*, and to deactivate them when we switch to *unicode font mode*. This is made (almost) entirely in lua. The difficult part is the changes in `\DeclareFontEncoding`.

```

366
367 \ProvidesFile{lutf8x.def}
368     [2010/05/10 v0.97 UTF-8 support for luainputenc]
369

```

```

370 \makeatletter
371 \catcode`\'\ \saved@space@catcode
372
373 \@inpcnt@test
374
375 \ifx\@begindocumenthook\@undefined
376   \makeatother
377   \endinput \fi
378

```

We change it a little to add the activated character in the lua table.

```

379
380 \gdef\DeclareUnicodeCharacter#1#2{%
381   \tempcnta"#1%
382   \luadirect{\luainputenc.declare_character('`the`\tempcnta')}%
383   \catcode`\tempcnta\active %
384   \DeclareInputText{\the\tempcnta}{#2}%
385 }
386
387 \onlypreamble\DeclareUnicodeCharacter
388
389 \def\cdp@elt#1#2#3#4{%
390   \wlog{Now handling font encoding #1 ...}%
391   \lowercase{%
392     \InputIfFileExists{#1enc.dfu}}%
393     {\wlog{... processing UTF-8 mapping file for font encoding
394           #1}%
395     \catcode`\ 9\relax}%
396     {\wlog{... no UTF-8 mapping file for font encoding #1}}%
397 }
398 \cdp@list
399

```

The macros to change from/to *legacy mode* to/from *unicode font mode*.

```

400
401 \def\lIE@ActivateUnicodeCatcodes{%
402 \luadirect{\luainputenc.activate_characters()}%
403 }
404
405 \def\lIE@DesactivateUnicodeCatcodes{%
406 \luadirect{\luainputenc.desactivate_characters()}%
407 }
408
409 \def\lIE@CharactersActivated{%
410 \luadirect{\luainputenc.force_characters_activated()}%
411 }
412
413 \edef\lIE@EU{EU2}
414

```

We add some code to automatically activate or deactivate characters according to the encoding changes. Note that we override `\@enc@update`, which may pose some problems

if a package of yours does it too. Fortunately this package is the only one that does it in TeXLive.

```

415
416 \def\DeclareFontEncoding@#1#2#3{%
417   \edef\lIE@test{#1}%
418   \ifx\lIE@test\lIE@EU %
419     \ifx\LastDeclaredEncoding\lIE@EU\else %
420       \lIE@CharactersActivated %
421       \lIE@DesactivateUnicodeCatcodes %
422     \fi
423   \gdef\@enc@update{%
424     \edef\lIE@test{#1}%
425     \ifx\f@encoding\lIE@EU %
426       \lIE@DesactivateUnicodeCatcodes %
427     \else %
428       \lIE@ActivateUnicodeCatcodes %
429     \fi
430     \expandafter\let\csname cf@encoding-cmd\endcsname@\changed@cmd
431     \expandafter\let\csname f@encoding-cmd\endcsname@\current@cmd
432     \default@T
433     \csname T@\f@encoding\endcsname
434     \csname D@\f@encoding\endcsname
435     \let\enc@update\relax
436     \let\cf@encoding\f@encoding
437   }
438 \else %
439   \expandafter %
440   \ifx\csname T@#1\endcsname\relax %
441     \def\cdp@elt{\noexpand\cdp@elt}%
442     \xdef\cdp@list{\cdp@list\cdp@elt{#1}%
443       {\default@family}{\default@series}%
444       {\default@shape}}%
445     \expandafter\let\csname#1-cmd\endcsname@\changed@cmd %
446     \begingroup %
447       \wlog{Now handling font encoding #1 ...}%
448       \lowercase{%
449         \InputIfFileExists{\enc.dfu}%
450           {\wlog{... processing UTF-8 mapping file for font encoding #1}}%
451           {\wlog{... no UTF-8 mapping file for font encoding #1}}%
452       \endgroup
453     \else
454       \font@info{Redeclaring font encoding #1}%
455     \fi
456   \fi %
457   \global\@namedef{T@#1}{#2}%
458   \global\@namedef{M@#1}{\default@M#3}%
459   \xdef\LastDeclaredEncoding{#1}%
460 }
461
462 \DeclareUnicodeCharacter{00A9}{\textcopyright}

```

```

463 \DeclareUnicodeCharacter{00AA}{\textordfeminine}
464 \DeclareUnicodeCharacter{00AE}{\textregistered}
465 \DeclareUnicodeCharacter{00BA}{\textordmasculine}
466 \DeclareUnicodeCharacter{02C6}{\textasciicircum}
467 \DeclareUnicodeCharacter{02DC}{\textasciitilde}
468 \DeclareUnicodeCharacter{200C}{\textcompwordmark}
469 \DeclareUnicodeCharacter{2026}{\textellipsis}
470 \DeclareUnicodeCharacter{2122}{\texttrademark}
471 \DeclareUnicodeCharacter{2423}{\textvisiblespace}
472

```

4.5 luainputenc.lua

First the `inputenc` module is registered as a LuaTeX module, with some informations.

```

473
474 module('luainputenc', package.seeall)
475
476 luainputenc.module = {
477     name      = "luainputenc",
478     version   = 0.97,
479     date      = "2010/05/10",
480     description = "Lua simple inputenc package.",
481     author    = "Elie Roux",
482     copyright = "Elie Roux",
483     license   = "CC0",
484 }
485
486 luatexbase.provides_module(luainputenc.module)
487
488 local format = string.format
489
490 luainputenc.log = luainputenc.log or function(...)
491     luatexbase.module_log('luainputenc', format(...))
492 end
493

```

We keep the option and the true encoding in two variables.

```

494
495 luainputenc.encoding = "utf8"
496 luainputenc.package_option = nil
497
498 function luainputenc.set_option(option)
499     luainputenc.package_option = option
500     if option == "lutf8" or option == "lutf8x" or option == "utf8x" or option == "unactivate" then
501         luainputenc.encoding = "utf8"
502     else
503         luainputenc.encoding = option
504     end
505 end
506

```

Some local declarations.

```
507
508 local char, utfchar, byte, format, gsub, utfbyte, utfgsub =
509 string.char, unicode.utf8.char, string.byte, string.format, string.gsub, unicode.utf8.byte, unicode.utf8.sub
```

510

The function to transform a 8-bit character in the corresponding fake UTF-8 character.

```
511
512 function luainputenc.byte_to_utf(ch)
513     return utfchar(byte(ch))
514 end
515
```

The function that will be registered in the `process_input_buffer` callback when needed.

```
516
517 function luainputenc.fake_utf_read(buf)
518     return gsub(buf,"(.)", luainputenc.byte_to_utf)
519 end
520
```

The function to transform a fake utf8 character in the corresponding 8-bit character.

```
521
522 function luainputenc.utf_to_byte(ch)
523     return char(utfbyte(ch))
524 end
525
```

The function that will be registered in the `process_output_buffer` callback if it exists.

```
526
527 function luainputenc.fake_utf_write(buf)
528     return utfgsub(buf,"(.)", luainputenc.utf_to_byte)
529 end
530
```

Here we register the two callbacks, and the behaviour is the same as in pdfTeX. The next part of the file is only the machinery for LuaTeX versions without the callback `process_output_buffer`, so it will be deprecated after TeXLive 2009, you are not advised to use it.

```
531
532 if tex.luatexversion > 42 then
533
534     function luainputenc.register_callbacks()
535         luatexbase.add_to_callback('process_output_buffer', luainputenc.fake_utf_write, 'luainputenc.register_callbacks')
536         luatexbase.add_to_callback('process_input_buffer', luainputenc.fake_utf_read, 'luainputenc.register_callbacks')
537     end
538
539 else
540
```

`start()` and `stop()` are the functions that register or unregister the function in the callback. When the function is registered, LuaTeX reads the input in fake UTF-8.

```

541
542     local started, stopped = 1, 0
543
544     luainputenc.state = stopped
545
546     function luainputenc.setstate(state)
547         if state == luainputenc.state then
548             return
549         elseif state == started then
550             luainputenc.start()
551         else
552             luainputenc.stop()
553         end
554     end
555
556     function luainputenc.setstarted()
557         luainputenc.setstate(started)
558     end
559
560     function luainputenc.setstopped()
561         luainputenc.setstate(stopped)
562     end
563
564     function luainputenc.start()
565         luatexbase.add_to_callback('process_input_buffer', luainputenc.fake_utf_read,
566             'luainputenc.fake_utf_read')
567         luainputenc.state = started
568         if luainputenc.callback_registered == 0 then
569             luainputenc.register_callback()
570         end
571     end
572
573     function luainputenc.stop()
574         luatexbase.remove_from_callback('process_input_buffer', 'luainputenc.fake_utf_read')
575         luainputenc.state = stopped
576         return
577     end
578

```

Here is a list of all file extention for which we consider that the files have been written by LuaTeX, and thus must be read in fake UTF-8. I may have forgotten things in the list. If you find a new extention, please report the maintainer.

```

579
580     luainputenc_unicode_extentions = {
581         ['.aux'] = 1, -- basic files
582         ['.toc'] = 1,
583         ['.gls'] = 1,
584         ['.ind'] = 1,

```

```

585     ['.idx'] = 1,
586     ['.vrb'] = 1, -- beamer and powerdot
587     ['.nav'] = 1, -- other beamer extention
588     ['.sol'] = 1,
589     ['.qsl'] = 1,
590     ['.snm'] = 1,
591     ['.pgn'] = 1, -- pagerefERENCE
592     ['.cpG'] = 1, -- AlProTeX
593     ['.pst'] = 1, -- pst-tree
594     ['.tmp'] = 1, -- sauerj/collect
595     ['.sym'] = 1, -- listofsymbols
596     ['.sub'] = 1, -- listofsymbols
597     ['.lof'] = 1, -- preprint
598     ['.lot'] = 1, -- preprint
599     ['.mtc1'] = 1, -- minitoc
600     ['.ovr'] = 1, -- thumbss
601     ['.fff'] = 1, -- endplate
602     ['.sbb'] = 1, -- splitbib
603     ['.bbl'] = 1, -- latex
604     ['.ain'] = 1, -- authorindex
605     ['.abb'] = 1, -- juraabbrev
606     ['.ent'] = 1, -- endnotes
607     ['.end'] = 1, -- fn2end
608     ['.thm'] = 1, -- ntheorem
609     ['.xtr'] = 1, -- extract
610     ['.han'] = 1, -- linguhO
611     ['.bnd'] = 1, -- bibref
612     ['.bbl'] = 1, -- bibref
613     ['.col'] = 1, -- mwrite
614     ['.ttt'] = 1, -- endfloat
615     ['.fax'] = 1, -- lettre
616     ['.tns'] = 1, -- lettre
617     ['.odt'] = 1, -- lettre
618     ['.etq'] = 1, -- lettre
619     ['.emd'] = 1, -- poemscol
620     ['.emx'] = 1, -- poemscol
621     ['.ctn'] = 1, -- poemscol
622     ['.hst'] = 1, -- vhISTORY
623     ['.acr'] = 1, -- crosswrd
624     ['.dwn'] = 1, -- crosswrd
625     ['.ttc'] = 1, -- talk
626     -- ['.txt'] = 1, -- coverpage, but not sure it's safe to include it...
627     ['.eve'] = 1, -- calend0
628     ['.scn'] = 1, -- cwebmac
629   }
630

```

The code to define a specific behaviour for certain files.

```

631
632     luainputenc_unicode_files = {}
633

```

```

634     luainputenc.non_unicode_files = {}
635
636     function luainputenc.set_unicode_file(filename)
637         if luainputenc.non_unicode_files[filename] == 1 then
638             luainputenc.non_unicode_files[filename] = nil
639         end
640         luainputenc_unicode_files[filename] = 1
641     end
642
643     function luainputenc.set_non_unicode_file(filename)
644         if luainputenc_unicode_files[filename] == 1 then
645             luainputenc_unicode_files[filename] = nil
646         end
647         luainputenc.non_unicode_files[filename] = 1
648     end
649
650     function luainputenc.set_unicode_extention(ext)
651         luainputenc_unicode_extention[ext] = 1
652     end
653
654     function luainputenc.set_non_unicode_extention(ext)
655         if luainputenc_unicode_extentions[ext] == 1 then
656             luainputenc_unicode_extentions[ext] = nil
657         end
658     end
659
660     function luainputenc.unset_file(filename)
661         if luainputenc_unicode_files[filename] == 1 then
662             luainputenc_unicode_files[filename] = nil
663         elseif luainputenc.non_unicode_files[filename] == 1 then
664             luainputenc.non_unicode_files[filename] = nil
665         end
666     end
667
668     local unicode, non_unicode = stopped, started
669
670     function luainputenc.find_state(filename)
671         if luainputenc_unicode_files[filename] == 1 then
672             return unicode
673         elseif luainputenc.non_unicode_files[filename] == 1 then
674             return non_unicode
675         else
676             local ext = filename:sub(-4)
677             if luainputenc_unicode_extentions[ext] == 1 then
678                 return unicode
679             else
680                 return non_unicode
681             end
682         end
683     end

```

684

We register the functions to stop or start the fake UTF-8 translation in the appropriate callbacks if necessary.

```
685     function luainputenc.pre_read_file(env)
686         if not env.path then
687             return
688         end
689         local currentstate = luainputenc.state
690         luainputenc.setstate(luainputenc.find_state(env.filename))
691         env.previousstate = currentstate
692     end
693
694
695     function luainputenc.close(env)
696         luainputenc.setstate(env.previousstate)
697     end
698
699     luainputenc.callback_registered = 0
700
701     function luainputenc.register_callback()
702         if luainputenc.callback_registered == 0 then
703             luatexbase.add_to_callback('pre_read_file', luainputenc.pre_read_file,
704                                         'luainputenc.pre_read_file')
705             luatexbase.add_to_callback('file_close', luainputenc.close, 'luainputenc.close')
706             luainputenc.callback_registered = 1
707         end
708     end
709
710 end
711
```

Finally we provide some functions to activate or deactivate the catcodes of the non-ASCII characters.

```
712
713
714 luainputenc.activated_characters = {}
715 luainputenc.characters_are_activated = false
716
717 function luainputenc.declare_character(c)
718     luainputenc.activated_characters[tonumber(c)] = true
719 end
720
721 function luainputenc.force_characters_activated ()
722     luainputenc.characters_are_activated = true
723 end
724
725 function luainputenc.activate_characters()
726     if not luainputenc.characters_are_activated then
727         for n, _ in pairs(luainputenc.activated_characters) do
```

```

728         tex.sprint(string.format('`\catcode %d\active',n))
729     end
730     luainputenc.characters_are_activated = true
731   end
732 end
733
734 function luainputenc.desactivate_characters()
735   if luainputenc.characters_are_activated then
736     for n, _ in pairs(luainputenc.activated_characters) do
737       tex.sprint(string.format('`\catcode %d=11',n))
738     end
739     luainputenc.characters_are_activated = false
740   end
741 end
742

```

5 Test file

Very minimal, just check that the package correctly loads.

```

743 <*test>
744 \RequirePackage{luainputenc}
745 \stop
746 </test>

```