

# The `luainputenc` package

Manuel Pégourié-Gonnard [mpg@elzevir.fr](mailto:mpg@elzevir.fr)  
Élie Roux [elie.roux@telecom-bretagne.eu](mailto:elie.roux@telecom-bretagne.eu)

2010/11/19 v0.973

## 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

<b>1 Overview: When (not) to use this package</b>	<b>1</b>
<b>2 Documentation</b>	<b>2</b>
2.1 Introduction . . . . .	2
2.2 Overview of 8-bit mode . . . . .	3
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
<b>3 Accessing the encoding in lua</b>	<b>4</b>
<b>4 Files</b>	<b>4</b>
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
<b>5 Test file</b>	<b>22</b>

## 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<sup>A</sup>T<sub>E</sub>X implements font encodings.

From the user point of view, adapting an old document for LuaT<sub>E</sub>X 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<sub>E</sub>X without changing them.

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

Unicode fonts with LuaT<sub>E</sub>X 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<sup>A</sup>T<sub>E</sub>X'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<sub>E</sub>X 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<sub>E</sub>X) do, and thus `inputenc` is totally broken with LuaT<sub>E</sub>X. This package aims at replacing `inputenc` for LuaT<sub>E</sub>X, by adapting the way LuaT<sub>E</sub>X 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 LuaTeX 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 (L<sup>A</sup>T<sub>E</sub>X Internal Character Representation), with the font encoding.

In LuaTeX version 0.43, a new callback called `process_output_buffer`, this callbacks allows to make LuaTeX 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 LuaTeX 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 LuaTeX 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, LuaTeX 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 LuaTeX 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 LuaTeX can read more than 256 characters is that fonts can also have more than 256 characters. LuaTeX 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 `luafontload`.

### 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<sup>A</sup>T<sub>E</sub>X 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<sup>A</sup>T<sub>E</sub>X 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` and ensure a few primitives, then we load the `lua` file.

```

67
68 \RequirePackage{luatexbase}
69 \luatexbase@ensure@primitive{luaescapestring}
70
71 \RequireLuaModule{luainputenc}
72

```

Here is some code from `inputenc`.

```

73
74 \def\DeclareInputMath#1{%
75   \@inpenc@test
76   \bgroup
77     \uccode`~#1%
78     \uppercase{%
79       \egroup
80       \def~{%
81         }%
82     }%
83 \def\DeclareInputText#1#2{%
84   \def\reserved@a##1 ${}%
85   \def\reserved@b{#2}%
86   \ifcat_ \expandafter\reserved@a\meaning\reserved@b$ ${}%
87     \DeclareInputMath{#1}{#2}%
88   \else
89     \DeclareInputMath{#1}{\IeC{#2}}%

```

```

90     \fi
91 }
92 \def\IeC{%
93   \ifx\protect\@typeset@protect
94     \expandafter\@firstofone
95   \else
96     \noexpand\IeC
97   \fi
98 }

```

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.

```

99
100 \def\inputencoding#1{%
101   \the\inpenc@prehook
102   \gdef\@inpenc@test{\global\let\@inpenc@test\relax}%
103   \edef\@inpenc@undefined{\noexpand\@inpenc@undefined@{#1}}%
104   \edef\inputencodingname{#1}%
105   \@inpenc@loop\^\^\A\^\^H%
106   \@inpenc@loop\^\^\K\^\^K%
107   \@inpenc@loop\^\^\N\^\^_%
108   \xdef\@saved@endlinechar{\the\endlinechar }%
109   \endlinechar=-1
110   \xdef\@saved@space@catcode{\the\catcode`\ }%
111   \catcode`\ 9\relax
112   \input{#1.def}%
113   \endlinechar=\@saved@endlinechar{}%
114   \catcode`\ \@saved@space@catcode\relax
115   \ifx\@inpenc@test\relax\else
116     \PackageWarning{inputenc}%
117       {No characters defined\MessageBreak
118        by input encoding change to '#1'\MessageBreak}%
119   \fi
120   \the\inpenc@posthook
121   \luatexbase@directlua{luainputenc.set_option("\luatexluaescapestring{#1}")}%
122 }
123 \newtoks\inpenc@prehook
124 \newtoks\inpenc@posthook
125 \def\@inpenc@undefined@#1{\PackageError{inputenc}%
126   {Keyboard character used is undefined\MessageBreak
127    in inputencoding '#1'}%
128   {You need to provide a definition with
129    \noexpand\DeclareInputText\MessageBreak or
130    \noexpand\DeclareInputMath before using this key.}%
131 \def\@inpenc@loop#1#2{%
132   \tempcnta`#1\relax

```

```

133 \loop
134   \catcode`\@tempcnta\active
135   \bgroup
136     \uccode`~\@tempcnta
137     \uppercase{%
138       \egroup
139       \let~\inpc@undefined
140     }%
141   \ifnum\@tempcnta<#2\relax
142     \advance\@tempcnta@ne
143   \repeat
144

```

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

```

145
146 \DeclareOption{utf8}{%
147   \inputencoding{lutf8}%
148 }
149
150 \DeclareOption{lutf8}{%
151   \inputencoding{lutf8}%
152 }
153
154 \DeclareOption{utf8x}{%
155   \inputencoding{lutf8}%
156 }
157
158 \DeclareOption{lutf8x}{%
159   \inputencoding{lutf8x}%
160 }
161

```

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

```

162
163 \DeclareOption{unactivate}{%
164   \edef\inputencodingname{unactivate}%
165   \luatexbase@directlua{luainputenc.set_option([[unactivate]])}
166 }
167

```

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.

```

168
169 \DeclareOption*{%
170   \lIE@activate %
171   \inpc@loop\^\^\?^\^\ff%
172   \inputencoding{\CurrentOption}%
173 }
174

```

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.

```

175
176 \ifnum\luatexversion>42
177
178   \newcommand*\lIE@activate[0]{%
179     \luatexbase@directlua{luainputenc.register_callbacks()}%
180   }
181
182 \else
183

```

`\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.

```

184
185 \newcommand*\lIE@setstarted[0]{%
186   \ifnum\lIE@activated=1 %
187     \luatexbase@directlua{luainputenc.setstarted()}%
188   \fi %
189 }
190
191 \newcommand*\lIE@setstopped[0]{%
192   \ifnum\lIE@activated=1 %
193     \luatexbase@directlua{luainputenc.setstopped()}%
194   \fi %
195 }
196

```

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.

```

197
198 \newcommand*\lIE@SetUtfFile[1]{%
199   \luatexbase@directlua{luainputenc.set_unicode_file("\luatexluaescapestring{\#1}")}%
200 }
201

```

\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).

```
202  
203 \newcommand*\LIE@SetNonUtfFile[1]{%  
204   \luatexbase@directlua{luainputenc.set_non_unicode_file("\luatexluaescapestring{\#1}")}%  
205 }  
206
```

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

```
207  
208 \newcommand*\LIE@UnsetFile[1]{%  
209   \luatexbase@directlua{luainputenc.unset_file("\luatexluaescapestring{\#1}")}%  
210 }  
211
```

\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.

```
212  
213 \newcommand*\LIE@SetUtfExt[1]{%  
214   \luatexbase@directlua{luainputenc.set_unicode_extention("\luatexluaescapestring{\#1}")}%  
215 }  
216
```

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

```
217  
218 \newcommand*\LIE@SetNonUtfExt[1]{  
219   \luatexbase@directlua{luainputenc.set_non_unicode_extention("\luatexluaescapestring{\#1}")}  
220 }  
221
```

\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.

```
222  
223  
224 \newcommand*\LIE@InputUtfFile[1]{%  
225   \LIE@SetUtfFile{\#1}%  
226   \input #1%  
227   \LIE@UnsetFile{\#1}%  
228 }  
229
```

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

```
230  
231 \newcommand*\LIE@InputNonUtfFile[1]{%  
232   \LIE@SetNonUtfFile{\#1}%
```

```

233 \input #1%
234 \lIE@UnsetFile{#1}%
235 }
236

```

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

```

237
238 \newcommand*\InputUtfFile[1]{%
239   \lIE@InputUtfFile{#1}%
240 }
241
242 \newcommand*\InputNonUtfFile[1]{%
243   \lIE@InputNonUtfFile{#1}%
244 }
245
246 \newcount\lIE@activated
247
248 \newcommand*{\lIE@activate}[0]{%
249   \lIE@activated=1 %
250   \lIE@setstarted %
251 }
252
253 \newcommand*{\lIE@FromInputenc}[1]{%
254   \ifnum\lIE@activated=0 %
255     \lIE@activate %
256   \fi%
257 }
258
259 \fi
260
261 \ProcessOptions*
262

```

### 4.3 lutf8.def

```

263 %% This file was adapted from utf8.def, which copyright is:
264 %% Copyright 1993 1994 1995 1996 1997 1998 1999 2000 2001 2002 2003
265 %% 2004 2005 2006 The LaTeX3 Project.
266 %%
267 %% utf8.def is under the lppl version 1.3c or later, and can be found
268 %% in the base LaTeX system.
269 %%
270 %% The lppl can be found at http://www.latex-project.org/lppl.txt
271 %%
272 %% The changes to utf8.def are Copyright 2009 Elie Roux, and are under
273 %% the CCO license.
274 %%
275 %% The changes are LuaTeX support.
276 %%
277 %% This file is distributed under the CCO license, with clause 6 of the

```

```

278 %% lppl as additional restrictions.
279

```

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.

```

280
281
282 \ProvidesFile{lutf8.def}
283   [2010/05/10 v0.97 UTF-8 support for luainputenc]
284
285 \makeatletter
286 \catcode`\\ \saved@space@catcode
287
288 \@inpcnt@test
289
290 \ifx\@begindocumenthook\@undefined
291   \makeatother
292   \endinput \fi
293

```

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`.

```

294
295 \gdef\DeclareUnicodeCharacter#1#2{%
296   \tempcnta"#1%
297   \catcode\tempcnta\active %
298   \DeclareInputText{\the\tempcnta}{#2}%
299 }
300
301 \onlypreamble\DeclareUnicodeCharacter
302
303 \def\cdp@elt#1#2#3#4{%
304   \wlog{Now handling font encoding #1 ...}%
305   \lowercase{%
306     \InputIfFileExists{#1enc.dfu}{}{%
307       \wlog{... processing UTF-8 mapping file for font encoding
308             #1}%
309       \catcode`\\ 9\relax}%
310     {\wlog{... no UTF-8 mapping file for font encoding #1}}%
311   }%
312 \cdp@list
313
314 \def\DeclareFontEncoding@#1#2#3{%
315   \expandafter %
316   \ifx\csname T#1\endcsname\relax %
317     \def\cdp@elt{\noexpand\cdp@elt}%
318     \xdef\cdp@list{\cdp@list\cdp@elt{#1}%
319                   {\default@family}{\default@series}%
320                   {\default@shape}}%

```

```

321 \expandafter\let\csname#1-cmd\endcsname\@changed@cmd %
322 \begingroup %
323   \wlog{Now handling font encoding #1 ...}%
324   \lowercase{%
325     \InputIfFileExists{#1enc.dfu}%
326       {\wlog{... processing UTF-8 mapping file for font encoding #1}}%
327       {\wlog{... no UTF-8 mapping file for font encoding #1}}%
328   \endgroup
329 \else
330   \font@info{Redeclaring font encoding #1}%
331 \fi
332 \global\@namedef{T@#1}{#2}%
333 \global\@namedef{M@#1}{\default@M#3}%
334 \xdef\LastDeclaredEncoding{#1}%
335 }
336
337 \DeclareUnicodeCharacter{00A9}{\textcopyright}
338 \DeclareUnicodeCharacter{00AA}{\textordfeminine}
339 \DeclareUnicodeCharacter{00AE}{\textregistered}
340 \DeclareUnicodeCharacter{00BA}{\textordmasculine}
341 \DeclareUnicodeCharacter{02C6}{\textasciicircum}
342 \DeclareUnicodeCharacter{02DC}{\textasciitilde}
343 \DeclareUnicodeCharacter{200C}{\textcompwordmark}
344 \DeclareUnicodeCharacter{2026}{\textellipsis}
345 \DeclareUnicodeCharacter{2122}{\texttrademark}
346 \DeclareUnicodeCharacter{2423}{\textvisiblespace}
347

```

#### 4.4 lutf8x.def

```

348 %% This file was adapted from utf8.def, which copyright is:
349 %% Copyright 1993 1994 1995 1996 1997 1998 1999 2000 2001 2002 2003
350 %% 2004 2005 2006 The LaTeX3 Project.
351 %%
352 %% utf8.def is under the lppl version 1.3c or later, and can be found
353 %% in the base LaTeX system.
354 %%
355 %% The lppl can be found at http://www.latex-project.org/lppl.txt
356 %%
357 %% The changes to utf8.def are Copyright 2009 Elie Roux, and are under
358 %% the CCO license.
359 %%
360 %% The changes are LuaTeX support.
361 %%
362 %% This file is distributed under the CCO license, with clause 6 of the
363 %% lppl as additional restrictions.
364

```

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 unactivate them when we switch to *unicode font mode*. This is

```

made (almost) entirely in lua. The difficult part is the changes in \DeclareFontEncoding.
365
366 \ProvidesFile{ltf8x.def}
367 [2010/05/10 v0.97 UTF-8 support for luainputenc]
368
369 \makeatletter
370 \catcode`\ \saved@space@catcode
371
372 \inpcntest
373
374 \ifx\@begindocumenthook\undefined
375   \makeatother
376   \endinput \fi
377

```

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

```

378
379 \gdef\DeclareUnicodeCharacter#1#2{%
380   \tempcnta"#1%
381   \luatexbase@directlua{luainputenc.declare_character('the\tempcnta')}%
382   \catcode\tempcnta\active %
383   \DeclareInputText{\the\tempcnta}{#2}%
384 }
385
386 \onlypreamble\DeclareUnicodeCharacter
387
388 \def\cdp@elt#1#2#3#4{%
389   \wlog{Now handling font encoding #1 ...}%
390   \lowercase{%
391     \InputIfFileExists{#1enc.dfu}%
392     {\wlog{... processing UTF-8 mapping file for font encoding
393       #1}%
394     \catcode`\ 9\relax}%
395     {\wlog{... no UTF-8 mapping file for font encoding #1}}%
396   }
397 \cdp@list
398

```

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

```

399
400 \def\lIE@ActivateUnicodeCatcodes{%
401 \luatexbase@directlua{luainputenc.activate_characters()}%
402 }
403
404 \def\lIE@DesactivateUnicodeCatcodes{%
405 \luatexbase@directlua{luainputenc.desactivate_characters()}%
406 }
407
408 \def\lIE@CharactersActivated{%
409 \luatexbase@directlua{luainputenc.force_characters_activated()}%
410 }

```

```

411
412 \edef\lIE@EU{EU2}
413
    We add some code to automatically activate or unactivate 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.
414
415 \def\DeclareFontEncoding#1#2#3{%
416   \edef\lIE@test{#1}%
417   \ifx\lIE@test\lIE@EU %
418     \ifx\LastDeclaredEncoding\lIE@EU\else %
419       \lIE@CharactersActivated %
420       \lIE@DesactivateUnicodeCatcodes %
421     \fi
422     \gdef\@@enc@update{%
423       \edef\lIE@test{#1}%
424       \ifx\f@encoding\lIE@EU %
425         \lIE@DesactivateUnicodeCatcodes %
426       \else %
427         \lIE@ActivateUnicodeCatcodes %
428       \fi
429       \expandafter\let\csname cf@encoding-cmd\endcsname\@changed@cmd
430       \expandafter\let\csname f@encoding-cmd\endcsname\@current@cmd
431       \default@T
432       \csname T@\f@encoding\endcsname
433       \csname D@\f@encoding\endcsname
434       \let\enc@update\relax
435       \let\cf@encoding\f@encoding
436     }
437   \else %
438     \expandafter %
439     \ifx\csname T@#1\endcsname\relax %
440       \def\cdp@elt{\noexpand\cdp@elt}%
441       \xdef\cdp@list{\cdp@list\cdp@elt{#1}%
442         {\default@family}{\default@series}%
443         {\default@shape}}%
444       \expandafter\let\csname#1-cmd\endcsname\@changed@cmd %
445       \begingroup %
446         \wlog{Now handling font encoding #1 ...}%
447         \lowercase{%
448           \InputIfFileExists{\enc.dfu}}%
449           {\wlog{... processing UTF-8 mapping file for font encoding #1}}%
450           {\wlog{... no UTF-8 mapping file for font encoding #1}}%
451         \endgroup
452       \else
453         \font@info{Redeclaring font encoding #1}%
454       \fi
455     \fi %
456   \global\@namedef{T@#1}{#2}%

```

```

457 \global\@namedef{M@#1}{\default@M#3}%
458 \xdef\LastDeclaredEncoding{#1}%
459 }
460
461 \DeclareUnicodeCharacter{00A9}{\textcopyright}
462 \DeclareUnicodeCharacter{00AA}{\textordfeminine}
463 \DeclareUnicodeCharacter{00AE}{\textregistered}
464 \DeclareUnicodeCharacter{00BA}{\textordmasculine}
465 \DeclareUnicodeCharacter{02C6}{\textasciicircum}
466 \DeclareUnicodeCharacter{02DC}{\textasciitilde}
467 \DeclareUnicodeCharacter{200C}{\textcompwordmark}
468 \DeclareUnicodeCharacter{2026}{\textellipsis}
469 \DeclareUnicodeCharacter{2122}{\texttrademark}
470 \DeclareUnicodeCharacter{2423}{\textvisiblespace}
471

```

## 4.5 luainputenc.lua

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

```

472
473 module('luainputenc', package.seeall)
474
475 luainputenc.module = {
476     name      = "luainputenc",
477     version   = 0.97,
478     date      = "2010/05/10",
479     description = "Lua simple inputenc package.",
480     author    = "Elie Roux",
481     copyright = "Elie Roux",
482     license   = "CC0",
483 }
484
485 luatexbase.provides_module(luainputenc.module)
486
487 local format = string.format
488
489 luainputenc.log = luainputenc.log or function(...)
490   luatexbase.module_log('luainputenc', format(...))
491 end
492

```

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

```

493
494 luainputenc.encoding = "utf8"
495 luainputenc.package_option = nil
496
497 function luainputenc.set_option(option)
498   luainputenc.package_option = option
499   if option == "lutf8" or option == "lutf8x" or option == "utf8x" or option == "unactivate" then
500     luainputenc.encoding = "utf8"
501   else

```

```
502     luainputenc.encoding = option
503   end
504 end
505
```

Some local declarations.

```
506
507 local char, utfchar, byte, format, gsub, utfbyte, utfgsub =
508 string.char, unicode.utf8.char, string.byte, string.format, string.gsub, unicode.utf8.byte, unicod
509
```

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

```
510
511 function luainputenc.byte_to_utf(ch)
512   return utfchar(byte(ch))
513 end
514
```

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

```
515
516 function luainputenc.fake_utf_read(buf)
517   return gsub(buf,"(.)", luainputenc.byte_to_utf)
518 end
519
```

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

```
520
521 function luainputenc.utf_to_byte(ch)
522   return char(utfbyte(ch))
523 end
524
```

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

```
525
526 function luainputenc.fake_utf_write(buf)
527   return utfgsub(buf,"(.)", luainputenc.utf_to_byte)
528 end
529
```

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.

```
530
531 if tex.luatexversion > 42 then
532
533   function luainputenc.register_callbacks()
534     luatexbase.add_to_callback('process_output_buffer', luainputenc.fake_utf_write, 'luainput
535     luatexbase.add_to_callback('process_input_buffer', luainputenc.fake_utf_read, 'luainputenc
536   end
```

```

537
538 else
539
      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.

540
541     local started, stopped = 1, 0
542
543     luainputenc.state = stopped
544
545     function luainputenc.setstate(state)
546         if state == luainputenc.state then
547             return
548         elseif state == started then
549             luainputenc.start()
550         else
551             luainputenc.stop()
552         end
553     end
554
555     function luainputenc.setstarted()
556         luainputenc.setstate(started)
557     end
558
559     function luainputenc.setstopped()
560         luainputenc.setstate(stopped)
561     end
562
563     function luainputenc.start()
564         luatexbase.add_to_callback('process_input_buffer', luainputenc.fake_utf_read,
565             'luainputenc.fake_utf_read')
566         luainputenc.state = started
567         if luainputenc.callback_registered == 0 then
568             luainputenc.register_callback()
569         end
570     end
571
572     function luainputenc.stop()
573         luatexbase.remove_from_callback('process_input_buffer', 'luainputenc.fake_utf_read')
574         luainputenc.state = stopped
575         return
576     end
577

```

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.

```

578
579     luainputenc_unicode_extensions = {

```

```

580     ['.aux'] = 1, -- basic files
581     ['.toc'] = 1,
582     ['.gls'] = 1,
583     ['.ind'] = 1,
584     ['.idx'] = 1,
585     ['.vrb'] = 1, -- beamer and powerdot
586     ['.nav'] = 1, -- other beamer extention
587     ['.sol'] = 1,
588     ['.qsl'] = 1,
589     ['.snm'] = 1,
590     ['.pgn'] = 1, -- pagereference
591     ['.cpg'] = 1, -- AlProTeX
592     ['.pst'] = 1, -- pst-tree
593     ['.tmp'] = 1, -- sauverj/collect
594     ['.sym'] = 1, -- listoftsymbols
595     ['.sub'] = 1, -- listoftsymbols
596     ['.lof'] = 1, -- preprint
597     ['.lot'] = 1, -- preprint
598     ['.mtc1'] = 1, -- minitoc
599     ['.ovr'] = 1, -- thumbss
600     ['.fff'] = 1, -- endplate
601     ['.sbb'] = 1, -- splitbib
602     ['.bb1'] = 1, -- latex
603     ['.ain'] = 1, -- authorindex
604     ['.abb'] = 1, -- juraabbrev
605     ['.ent'] = 1, -- endnotes
606     ['.end'] = 1, -- fn2end
607     ['.thm'] = 1, -- ntheorem
608     ['.xtr'] = 1, -- extract
609     ['.han'] = 1, -- linguhoh
610     ['.bnd'] = 1, -- bibref
611     ['.bb1'] = 1, -- bibref
612     ['.col'] = 1, -- mwwrite
613     ['.ttt'] = 1, -- endfloat
614     ['.fax'] = 1, -- lettre
615     ['.tns'] = 1, -- lettre
616     ['.odt'] = 1, -- lettre
617     ['.etq'] = 1, -- lettre
618     ['.emd'] = 1, -- poemscol
619     ['.emx'] = 1, -- poemscol
620     ['.ctn'] = 1, -- poemscol
621     ['.hst'] = 1, -- vhhistory
622     ['.acr'] = 1, -- crosswrd
623     ['.dwn'] = 1, -- crosswrd
624     ['.ttc'] = 1, -- talk
625     -- ['.txt'] = 1, -- coverpage, but not sure it's safe to include it...
626     ['.eve'] = 1, -- calend0
627     ['.scn'] = 1, -- cwebmac
628 }
629

```

The code to define a specific behaviour for certain files.

```
630
631     luainputenc_unicode_files = {}
632
633     luainputenc_non_unicode_files = {}
634
635     function luainputenc.set_unicode_file(filename)
636         if luainputenc.non_unicode_files[filename] == 1 then
637             luainputenc.non_unicode_files[filename] = nil
638         end
639         luainputenc_unicode_files[filename] = 1
640     end
641
642     function luainputenc.set_non_unicode_file(filename)
643         if luainputenc_unicode_files[filename] == 1 then
644             luainputenc_unicode_files[filename] = nil
645         end
646         luainputenc_non_unicode_files[filename] = 1
647     end
648
649     function luainputenc.set_unicode_extention(ext)
650         luainputenc_unicode_extention[ext] = 1
651     end
652
653     function luainputenc.set_non_unicode_extention(ext)
654         if luainputenc_unicode_extentions[ext] == 1 then
655             luainputenc_unicode_extentions[ext] = nil
656         end
657     end
658
659     function luainputenc.unset_file(filename)
660         if luainputenc_unicode_files[filename] == 1 then
661             luainputenc_unicode_files[filename] = nil
662         elseif luainputenc_non_unicode_files[filename] == 1 then
663             luainputenc_non_unicode_files[filename] = nil
664         end
665     end
666
667     local unicode, non_unicode = stopped, started
668
669     function luainputenc.find_state(filename)
670         if luainputenc_unicode_files[filename] == 1 then
671             return unicode
672         elseif luainputenc_non_unicode_files[filename] == 1 then
673             return non_unicode
674         else
675             local ext = filename:sub(-4)
676             if luainputenc_unicode_extentions[ext] == 1 then
677                 return unicode
678             else
```

```

679             return non_unicode
680         end
681     end
682 end
683

```

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

```

684
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     function luainputenc.close(env)
695         luainputenc.setstate(env.previousstate)
696     end
697
698     luainputenc.callback_registered = 0
699
700     function luainputenc.register_callback()
701         if luainputenc.callback_registered == 0 then
702             luatexbase.add_to_callback('pre_read_file', luainputenc.pre_read_file,
703             'luainputenc.pre_read_file')
704             luatexbase.add_to_callback('file_close', luainputenc.close, 'luainputenc.close')
705             luainputenc.callback_registered = 1
706         end
707     end
708
709 end
710

```

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

```

711
712
713 luainputenc.activated_characters = {}
714 luainputenc.characters_are_activated = false
715
716 function luainputenc.declare_character(c)
717     luainputenc.activated_characters[tonumber(c)] = true
718 end
719
720 function luainputenc.force_characters_activated ()
721     luainputenc.characters_are_activated = true
722 end

```

```

723
724 function luainputenc.activate_characters()
725     if not luainputenc.characters_are_activated then
726         for n, _ in pairs(luainputenc.activated_characters) do
727             tex.sprint(string.format('\\catcode %d\\active',n))
728         end
729         luainputenc.characters_are_activated = true
730     end
731 end
732
733 function luainputenc.desactivate_characters()
734     if luainputenc.characters_are_activated then
735         for n, _ in pairs(luainputenc.activated_characters) do
736             tex.sprint(string.format('\\catcode %d=11',n))
737         end
738         luainputenc.characters_are_activated = false
739     end
740 end
741

```

## 5 Test file

Very minimal, just check that the package correctly loads with an option and doesn't crash on a one-line plain ASCII document body...

```

742 (*test)
743 \documentclass{article}
744 \usepackage[utf8]{luainputenc}
745 \begin{document}
746 bla
747 \end{document}
748 
```