# The luainputenc package

Elie Roux

`elie.roux@telecom-bretagne.eu`

2009/09/23 v0.94

**Abstract**

Input encoding management for LuaTeX. For an introduction on this package (among others), please refer to `luatex-reference.pdf`.

## Contents

## 1 Documentation

### 1.1 Introduction

One the the most interesting new features of LuaTeX 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 pdfTeX) do, and thus inputenc is totally broken with LuaTeX. This package aims at replacing inputenc for LuaTeX, by adapting the way LuaTeX handles input, and the way inputenc handles UTF-8. This package has two very distinct modes: 8-bit and UTF-8.

## 1.2   Overview of 8-bit mode

This package **does not** map 8-bit encodings to utf8. It allows LuaTEX to read 8-bit charac-
ters, 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 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 extentions. 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.

## 1.3   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.).

### 1.3.1   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`.

### 1.3.2 unicode font mode

To fix this, the most simple solution is to desactivate 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 this package too. See section 2.5 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.

### 1.3.3 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.

## 2 Files

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

### 2.1 `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 LaTeX 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
```

## 2.2  `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 CC0 license.
37 %%
38 %% The changes are LuaTeX support.
39 %%
40 %% This file is distributed under the CC0 license, with clause 6 of the
41 %% lppl as additional restrictions.
42
```

First we check if we are called with LuaTEX, (pdf)TEXor 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   \RequirePackage{xetex-inputenc}
49   \DeclareOption{unactivate}{\PassOptionsToPackage{utf8}{xetex-inputenc}}
50   \DeclareOption{lutf8}{\PassOptionsToPackage{utf8}{xetex-inputenc}}
51   \DeclareOption{lutf8x}{\PassOptionsToPackage{utf8}{xetex-inputenc}}
52   \DeclareOption*{\PassOptionsToPackage{\CurrentOption}{xetex-inputenc}}
```

```
53   \ProcessOptions*
54   \expandafter\endinput
55 \fi
56
57 \ifluatex\else
58   \RequirePackage{inputenc}
59   \DeclareOption{unactivate}{\PassOptionsToPackage{utf8}{inputenc}}
60   \DeclareOption{lutf8}{\PassOptionsToPackage{utf8}{inputenc}}
61   \DeclareOption{lutf8x}{\PassOptionsToPackage{utf8}{inputenc}}
62   \DeclareOption*{\PassOptionsToPackage{\CurrentOption}{inputenc}}
63   \ProcessOptions*
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{luatextra}
69
70 \luatexUseModule{luainputenc}
71
```

Here is some code from inputenc.

```
72
73 \def\DeclareInputMath#1{%
74     \@inpenc@test
75     \bgroup
76         \uccode'\~#1%
77         \uppercase{%
78     \egroup
79         \def~%
80     }%
81 }
82 \def\DeclareInputText#1#2{%
83     \def\reserved@a##1 ${}%
84     \def\reserved@b{#2}%
85     \ifcat_\expandafter\reserved@a\meaning\reserved@b$ $_%
86         \DeclareInputMath{#1}{#2}%
87     \else
88         \DeclareInputMath{#1}{\IeC{#2}}%
89     \fi
90 }
91 \def\IeC{%
92   \ifx\protect\@typeset@protect
93     \expandafter\@firstofone
94   \else
95     \noexpand\IeC
96   \fi
97 }
```

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.

```
98
99 \def\inputencoding#1{%
100   \the\inpenc@prehook
101   \gdef\@inpenc@test{\global\let\@inpenc@test\relax}%
102   \edef\@inpenc@undefined{\noexpand\@inpenc@undefined@{#1}}%
103   \edef\inputencodingname{#1}%
104   \@inpenc@loop\^^A\^^H%
105   \@inpenc@loop\^^K\^^K%
106   \@inpenc@loop\^^N\^^_%
107   \advance\endlinechar\@M
108   \xdef\saved@space@catcode{\the\catcode`\ }%
109   \catcode`\ 9\relax
110   \input{#1.def}%
111   \advance\endlinechar-\@M
112   \catcode`\ \saved@space@catcode\relax
113   \ifx\@inpenc@test\relax\else
114     \PackageWarning{inputenc}%
115             {No characters defined\MessageBreak
116              by input encoding change to `#1'\MessageBreak}%
117   \fi
118   \the\inpenc@posthook
119 }
120 \newtoks\inpenc@prehook
121 \newtoks\inpenc@posthook
122 \def\@inpenc@undefined@#1{\PackageError{inputenc}%
123         {Keyboard character used is undefined\MessageBreak
124          in inputencoding `#1'}%
125        {You need to provide a  definition with
126         \noexpand\DeclareInputText\MessageBreak or
127         \noexpand\DeclareInputMath before using this key.}}%
128 \def\@inpenc@loop#1#2{%
129   \@tempcnta`#1\relax
130   \loop
131     \catcode\@tempcnta\active
132     \bgroup
133       \uccode`\~\@tempcnta
134       \uppercase{%
135     \egroup
136         \let~\inpenc@undefined
137       }%
138   \ifnum\@tempcnta<`#2\relax
139     \advance\@tempcnta\@ne
140   \repeat}
141
```

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

```
142
```

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

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

```
159
160 \DeclareOption{unactivate}{%
161   \edef\inputencodingname{unactivate}%
162 }
163
```

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`.

```
164
165 \DeclareOption*{%
166   \lIE@activate %
167   \@inpenc@loop\^^?\^^ff%
168   \inputencoding{\CurrentOption}%
169 }
170
```

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.

```
171
172 \ifnum\luatexversion>42
173
174     \newcommand*{\lIE@activate}[0]{%
175       \luadirect{luainputenc.register_callbacks()}%
176     }
177
178 \else
179
```

`\lIE@setstarted` and `\lIE@setstopped` are called when the fake UTF-8 translation must be activated or desactivated. 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.

```
180
181 \newcommand*\lIE@setstarted[0]{%
182   \ifnum\lIE@activated=1 %
183     \luadirect{luainputenc.setstarted()}%
184   \fi %
185 }
186
187 \newcommand*\lIE@setstopped[0]{%
188   \ifnum\lIE@activated=1 %
189     \luadirect{luainputenc.setstopped()}%
190   \fi %
191 }
192
```

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 extentions 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.

```
193
194 \newcommand*\lIE@SetUtfFile[1]{%
195   \luadirect{luainputenc.set_unicode_file([[#1]])}%
196 }
197
```

\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).

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

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

```
203
204 \newcommand*\lIE@UnsetFile[1]{%
205   \luadirect{luainputenc.unset_file([[#1]])}%
206 }
207
```

\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

8

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.

```
208
209 \newcommand*\lIE@SetUtfExt[1]{%
210   \luadirect{luainputenc.set_unicode_extention([[#1]])}%
211 }
212
```

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

```
213
214 \newcommand*\lIE@SetNonUtfExt[1]{
215   \luadirect{luainputenc.set_non_unicode_extention([[#1]])}
216 }
217
```

\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.

```
218
219
220 \newcommand*\lIE@InputUtfFile[1]{%
221   \lIE@SetUtfFile{#1}%
222   \input #1%
223   \lIE@UnsetFile{#1}%
224 }
225
```

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

```
226
227 \newcommand*\lIE@InputNonUtfFile[1]{%
228   \lIE@SetNonUtfFile{#1}%
229   \input #1%
230   \lIE@UnsetFile{#1}%
231 }
232
```

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

```
233
234 \newcommand*\InputUtfFile[1]{%
235   \lIE@InputUtfFile{#1}%
236 }
237
238 \newcommand*\InputNonUtfFile[1]{%
239   \lIE@InputNonUtfFile{#1}%
240 }
241
242 \newcount\lIE@activated
243
```

```
244 \newcommand*{\lIE@activate}[0]{%
245   \lIE@activated=1 %
246   \lIE@setstarted %
247 }
248
249 \newcommand*{\lIE@FromInputenc}[1]{%
250   \ifnum\lIE@activated=0 %
251     \lIE@activate %
252   \fi%
253 }
254
255 \fi
256
257 \ProcessOptions*
258
```

## 2.3   `lutf8.def`

```
259 %% This file was adapted from utf8.def, which copyright is:
260 %% Copyright 1993 1994 1995 1996 1997 1998 1999 2000 2001 2002 2003
261 %% 2004 2005 2006 The LaTeX3 Project.
262 %%
263 %% utf8.def is under the lppl version 1.3c or later, and can be found
264 %% in the base LaTeX system.
265 %%
266 %% The lppl can be found at http://www.latex-project.org/lppl.txt
267 %%
268 %% The changes to utf8.def are Copyright 2009 Elie Roux, and are under
269 %% the CC0 license.
270 %%
271 %% The changes are LuaTeX support.
272 %%
273 %% This file is distributed under the CC0 license, with clause 6 of the
274 %% lppl as additional restrictions.
275
```

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.

```
276
277
278 \ProvidesFile{lutf8.def}
279   [2009/09/23 v0.94 UTF-8 support for luainputenc]
280
281 \makeatletter
282 \catcode`\ \saved@space@catcode
283
284 \@inpenc@test
285
286 \ifx\@begindocumenthook\@undefined
287   \makeatother
288   \endinput \fi
```

289

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.

```
290
291 \gdef\DeclareUnicodeCharacter#1#2{%
292  \@tempcnta"#1%
293  \catcode\@tempcnta\active %
294  \DeclareInputText{\the\@tempcnta}{#2}%
295 }
296
297 \@onlypreamble\DeclareUnicodeCharacter
298
299 \def\cdp@elt#1#2#3#4{%
300    \wlog{Now handling font encoding #1 ...}%
301    \lowercase{%
302        \InputIfFileExists{#1enc.dfu}}%
303          {\wlog{... processing UTF-8 mapping file for font encoding
304                #1}%
305            \catcode`\ 9\relax}%
306          {\wlog{... no UTF-8 mapping file for font encoding #1}}%
307 }
308 \cdp@list
309
310 \def\DeclareFontEncoding@#1#2#3{%
311    \expandafter %
312    \ifx\csname T@#1\endcsname\relax %
313      \def\cdp@elt{\noexpand\cdp@elt}%
314      \xdef\cdp@list{\cdp@list\cdp@elt{#1}%
315                    {\default@family}{\default@series}%
316                    {\default@shape}}%
317      \expandafter\let\csname#1-cmd\endcsname\@changed@cmd %
318      \begingroup %
319        \wlog{Now handling font encoding #1 ...}%
320        \lowercase{%
321          \InputIfFileExists{#1enc.dfu}}%
322            {\wlog{... processing UTF-8 mapping file for font encoding #1}}%
323            {\wlog{... no UTF-8 mapping file for font encoding #1}}%
324      \endgroup
325    \else
326      \@font@info{Redeclaring font encoding #1}%
327    \fi
328    \global\@namedef{T@#1}{#2}%
329    \global\@namedef{M@#1}{\default@M#3}%
330    \xdef\LastDeclaredEncoding{#1}%
331 }
332
333 \DeclareUnicodeCharacter{00A9}{\textcopyright}
334 \DeclareUnicodeCharacter{00AA}{\textordfeminine}
```

```
335 \DeclareUnicodeCharacter{00AE}{\textregistered}
336 \DeclareUnicodeCharacter{00BA}{\textordmasculine}
337 \DeclareUnicodeCharacter{02C6}{\textasciicircum}
338 \DeclareUnicodeCharacter{02DC}{\textasciitilde}
339 \DeclareUnicodeCharacter{200C}{\textcompwordmark}
340 \DeclareUnicodeCharacter{2026}{\textellipsis}
341 \DeclareUnicodeCharacter{2122}{\texttrademark}
342 \DeclareUnicodeCharacter{2423}{\textvisiblespace}
343
```

## 2.4   lutf8x.def

```
344 %% This file was adapted from utf8.def, which copyright is:
345 %% Copyright 1993 1994 1995 1996 1997 1998 1999 2000 2001 2002 2003
346 %% 2004 2005 2006 The LaTeX3 Project.
347 %%
348 %% utf8.def is under the lppl version 1.3c or later, and can be found
349 %% in the base LaTeX system.
350 %%
351 %% The lppl can be found at http://www.latex-project.org/lppl.txt
352 %%
353 %% The changes to utf8.def are Copyright 2009 Elie Roux, and are under
354 %% the CC0 license.
355 %%
356 %% The changes are LuaTeX support.
357 %%
358 %% This file is distributed under the CC0 license, with clause 6 of the
359 %% lppl as additional restrictions.
360
```

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`.

```
361
362 \ProvidesFile{lutf8x.def}
363    [2009/09/23 v0.94 UTF-8 support for luainputenc]
364
365 \makeatletter
366 \catcode`\ \saved@space@catcode
367
368 \@inpenc@test
369
370 \ifx\@begindocumenthook\@undefined
371   \makeatother
372   \endinput \fi
373
```

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

```
374
375 \gdef\DeclareUnicodeCharacter#1#2{%
376 \@tempcnta"#1%
```

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

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

```
395
396  \def\lIE@ActivateUnicodeCatcodes{%
397  \luadirect{luainputenc.activate_characters()}%
398  }
399
400  \def\lIE@DesactivateUnicodeCatcodes{%
401  \luadirect{luainputenc.desactivate_characters()}%
402  }
403
404  \def\lIE@CharactersActivated{%
405  \luadirect{luainputenc.force_characters_activated()}
406  }
407
408  \edef\lIE@EU{EU2}
409
```

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.

```
410
411  \def\DeclareFontEncoding@#1#2#3{%
412    \edef\lIE@test{#1}%
413    \ifx\lIE@test\lIE@EU %
414      \ifx\LastDeclaredEncoding\lIE@EU\else %
415        \lIE@CharactersActivated %
416        \lIE@DesactivateUnicodeCatcodes %
417      \fi
418      \gdef\@@enc@update{%
419        \edef\lIE@test{#1}%
420        \ifx\f@encoding\lIE@EU %
```

```
421        \lIE@DesactivateUnicodeCatcodes %
422       \else %
423        \lIE@ActivateUnicodeCatcodes %
424      \fi
425      \expandafter\let\csname\cf@encoding-cmd\endcsname\@changed@cmd
426      \expandafter\let\csname\f@encoding-cmd\endcsname\@current@cmd
427      \default@T
428      \csname T@\f@encoding\endcsname
429      \csname D@\f@encoding\endcsname
430      \let\enc@update\relax
431      \let\cf@encoding\f@encoding
432    }
433  \else %
434     \expandafter %
435     \ifx\csname T@#1\endcsname\relax %
436       \def\cdp@elt{\noexpand\cdp@elt}%
437       \xdef\cdp@list{\cdp@list\cdp@elt{#1}%
438                     {\default@family}{\default@series}%
439                     {\default@shape}}%
440       \expandafter\let\csname#1-cmd\endcsname\@changed@cmd %
441       \begingroup %
442         \wlog{Now handling font encoding #1 ...}%
443         \lowercase{%
444           \InputIfFileExists{#1enc.dfu}}%
445             {\wlog{... processing UTF-8 mapping file for font encoding #1}}%
446             {\wlog{... no UTF-8 mapping file for font encoding #1}}%
447       \endgroup
448     \else
449       \@font@info{Redeclaring font encoding #1}%
450     \fi
451  \fi %
452  \global\@namedef{T@#1}{#2}%
453  \global\@namedef{M@#1}{\default@M#3}%
454  \xdef\LastDeclaredEncoding{#1}%
455 }
456
457 \DeclareUnicodeCharacter{00A9}{\textcopyright}
458 \DeclareUnicodeCharacter{00AA}{\textordfeminine}
459 \DeclareUnicodeCharacter{00AE}{\textregistered}
460 \DeclareUnicodeCharacter{00BA}{\textordmasculine}
461 \DeclareUnicodeCharacter{02C6}{\textasciicircum}
462 \DeclareUnicodeCharacter{02DC}{\textasciitilde}
463 \DeclareUnicodeCharacter{200C}{\textcompwordmark}
464 \DeclareUnicodeCharacter{2026}{\textellipsis}
465 \DeclareUnicodeCharacter{2122}{\texttrademark}
466 \DeclareUnicodeCharacter{2423}{\textvisiblespace}
467
```

## 2.5  `eu2enc.def`

This file is extremely short. It just declares the encoding, with the default font. The default font here is lmr, which means that LATEX will read `eu2lmr.fd`. The problem is that all unicode fonts are OTF fonts, so `eu2lmr.fd` will call OTF fonts. Thus, to use `EU2`, you need to be able to read OTF fonts. The package luaotfload is a good choice to be able to do so.

```
468
469 \ProvidesFile{eu2enc.def}[2009/09/23 v0.1 a unicode font encoding for LuaTeX.]
470 \DeclareFontEncoding{EU2}{}{}
471 \DeclareErrorFont{EU2}{lmr}{m}{n}{10}
472 \DeclareFontSubstitution{EU2}{lmr}{m}{n}
473
```

## 2.6  `eu2lmr.fd`

This file simply describes the default (lmr) font of the `EU2` encoding. It loads the otf fonts with some default features enabled. This file may change, don't rely on it too much.

```
474
475 \ProvidesFile{eu2lmr.fd}
476     [2009/09/23 v0.2 Font defs for Latin Modern for LuaTeX's EU2 encoding]
477 \DeclareFontFamily{EU2}{lmr}{}
478 \DeclareFontShape{EU2}{lmr}{m}{n}%
479     {<-5.5>    "lmroman5-regular:+tlig;+tsub;+liga;+rlig;"
480   <5.5-6.5> "lmroman6-regular:+tlig;+tsub;+liga;+rlig;"
481     <6.5-7.5> "lmroman7-regular:+tlig;+tsub;+liga;+rlig;"
482   <7.5-8.5> "lmroman8-regular:+tlig;+tsub;+liga;+rlig;"
483     <8.5-9.5> "lmroman9-regular:+tlig;+tsub;+liga;+rlig;"
484   <9.5-11>  "lmroman10-regular:+tlig;+tsub;+liga;+rlig;"
485     <11-15>    "lmroman12-regular:+tlig;+tsub;+liga;+rlig;"
486     <15->      "lmroman17-regular:+tlig;+tsub;+liga;+rlig;"
487     }{}
488 \DeclareFontShape{EU2}{lmr}{m}{sl}%
489     {<-8.5>    "lmroman8-oblique:+tlig;+tsub;+liga;+rlig;"
490   <8.5-9.5> "lmroman9-oblique:+tlig;+tsub;+liga;+rlig;"
491     <9.5-11>  "lmroman10-oblique:+tlig;+tsub;+liga;+rlig;"
492   <11-15>   "lmroman12-oblique:+tlig;+tsub;+liga;+rlig;"
493     <15->      "lmroman17-oblique:+tlig;+tsub;+liga;+rlig;"
494     }{}
495 \DeclareFontShape{EU2}{lmr}{m}{it}%
496     {<-7.5>    "lmroman7-italic:+tlig;+tsub;+liga;+rlig;"
497     <7.5-8.5> "lmroman8-italic:+tlig;+tsub;+liga;+rlig;"
498   <8.5-9.5> "lmroman9-italic:+tlig;+tsub;+liga;+rlig;"
499     <9.5-11>  "lmroman10-italic:+tlig;+tsub;+liga;+rlig;"
500   <11->      "lmroman12-italic:+tlig;+tsub;+liga;+rlig;"
501     }{}
502 \DeclareFontShape{EU2}{lmr}{m}{sc}%
503     {<-> "lmroman10-capsregular:+tlig;+tsub;+liga;+rlig;"}{}
504 %
505 % Is this the right 'shape'?:
```

```
506 \DeclareFontShape{EU2}{lmr}{m}{scsl}%
507      {<-> "lmroman10-capsoblique:+tlig;+tsub;+liga;+rlig;"}{}
508 %%%%%% bold series
509 \DeclareFontShape{EU2}{lmr}{b}{n}
510      {<-> "lmroman10-demi:+tlig;+tsub;+liga;+rlig;"}{}
511 \DeclareFontShape{EU2}{lmr}{b}{sl}
512      {<-> "lmroman10-demioblique:+tlig;+tsub;+liga;+rlig;"}{}
513 %%%%%%% bold extended series
514 \DeclareFontShape{EU2}{lmr}{bx}{n}
515      {<-5.5>    "lmroman5-bold:+tlig;+tsub;+liga;+rlig;"
516   <5.5-6.5> "lmroman6-bold:+tlig;+tsub;+liga;+rlig;"
517      <6.5-7.5> "lmroman7-bold:+tlig;+tsub;+liga;+rlig;"
518   <7.5-8.5> "lmroman8-bold:+tlig;+tsub;+liga;+rlig;"
519      <8.5-9.5> "lmroman9-bold:+tlig;+tsub;+liga;+rlig;"
520   <9.5-11>  "lmroman10-bold:+tlig;+tsub;+liga;+rlig;"
521      <11->     "lmroman12-bold:+tlig;+tsub;+liga;+rlig;"
522      }{}
523 \DeclareFontShape{EU2}{lmr}{bx}{it}
524      {<-> "lmroman10-bolditalic:+tlig;+tsub;+liga;+rlig;"}{}
525 \DeclareFontShape{EU2}{lmr}{bx}{sl}
526      {<-> "lmroman10-boldoblique:+tlig;+tsub;+liga;+rlig;"}{}
527
```

## 2.7 `luainputenc.lua`

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

```
528
529 luainputenc = { }
530
531 luainputenc.module = {
532     name          = "luainputenc",
533     version       =  0.94,
534     date          = "2009/09/23",
535     description   = "Lua simple inputenc package.",
536     author        = "Elie Roux",
537     copyright     = "Elie Roux",
538     license       = "CC0",
539 }
540
541 luatextra.provides_module(luainputenc.module)
542
543 local format = string.format
544
545 luainputenc.log = luainputenc.log or function(...)
546   luatextra.module_log('luainputenc', format(...))
547 end
548
549 local char, utfchar, byte, format, gsub, utfbyte, utfgsub =
550 string.char, unicode.utf8.char, string.byte, string.format, string.gsub, unicode.utf8.byte, unicod
551
```

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

```
552
553 function luainputenc.byte_to_utf(ch)
554     return utfchar(byte(ch))
555 end
556
```

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

```
557
558 function luainputenc.fake_utf_read(buf)
559     return gsub(buf,"(.)", luainputenc.byte_to_utf)
560 end
561
```

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

```
562
563 function luainputenc.utf_to_byte(ch)
564     return char(utfbyte(ch))
565 end
566
```

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

```
567
568 function luainputenc.fake_utf_write(buf)
569     return utfgsub(buf,"(.)", luainputenc.utf_to_byte)
570 end
571
```

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.

```
572
573 if tex.luatexversion > 42 then
574
575     function luainputenc.register_callbacks()
576         callback.add('process_output_buffer', luainputenc.fake_utf_write, 'luainputenc.fake_utf_wr
577         callback.add('process_input_buffer', luainputenc.fake_utf_read, 'luainputenc.fake_utf_read
578     end
579
580 else
581
```

`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.

```
582
583     local started, stopped = 1, 0
584
585     luainputenc.state = stopped
586
```

```
587    function luainputenc.setstate(state)
588        if state == luainputenc.state then
589            return
590        elseif state == started then
591            luainputenc.start()
592        else
593            luainputenc.stop()
594        end
595    end
596
597    function luainputenc.setstarted()
598        luainputenc.setstate(started)
599    end
600
601    function luainputenc.setstopped()
602        luainputenc.setstate(stopped)
603    end
604
605    function luainputenc.start()
606        callback.add('process_input_buffer', luainputenc.fake_utf_read,
607            'luainputenc.fake_utf_read')
608        luainputenc.state = started
609        if luainputenc.callback_registered == 0 then
610            luainputenc.register_callback()
611        end
612    end
613
614    function luainputenc.stop()
615        callback.remove('process_input_buffer', 'luainputenc.fake_utf_read')
616        luainputenc.state = stopped
617        return
618    end
619
```

Here is a list of all file extentions 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.

```
620
621    luainputenc.unicode_extentions = {
622      ['.aux'] = 1, -- basic files
623      ['.toc'] = 1,
624      ['.gls'] = 1,
625      ['.ind'] = 1,
626      ['.idx'] = 1,
627      ['.vrb'] = 1, -- beamer and powerdot
628      ['.nav'] = 1, -- other beamer extentions
629      ['.sol'] = 1,
630      ['.qsl'] = 1,
631      ['.snm'] = 1,
632      ['.pgn'] = 1, -- pagereference
```

18

```
633        ['.cpg'] = 1, -- AlProTeX
634        ['.pst'] = 1, -- pst-tree
635        ['.tmp'] = 1, -- sauerj/collect
636        ['.sym'] = 1, -- listofsymbols
637        ['.sub'] = 1, -- listofsymbols
638        ['.lof'] = 1, -- preprint
639        ['.lot'] = 1, -- preprint
640        ['mtc1'] = 1, -- minitoc
641        ['.ovr'] = 1, -- thumbss
642        ['.fff'] = 1, -- endplate
643        ['.sbb'] = 1, -- splitbib
644        ['.bbl'] = 1, -- latex
645        ['.ain'] = 1, -- authorindex
646        ['.abb'] = 1, -- juraabbrev
647        ['.ent'] = 1, -- endnotes
648        ['.end'] = 1, -- fn2end
649        ['.thm'] = 1, -- ntheorem
650        ['.xtr'] = 1, -- extract
651        ['.han'] = 1, -- linguho
652        ['.bnd'] = 1, -- bibref
653        ['.bbl'] = 1, -- bibref
654        ['.col'] = 1, -- mwrite
655        ['.ttt'] = 1, -- endfloat
656        ['.fax'] = 1, -- lettre
657        ['.tns'] = 1, -- lettre
658        ['.odt'] = 1, -- lettre
659        ['.etq'] = 1, -- lettre
660        ['.emd'] = 1, -- poemscol
661        ['.emx'] = 1, -- poemscol
662        ['.ctn'] = 1, -- poemscol
663        ['.hst'] = 1, -- vhistory
664        ['.acr'] = 1, -- crosswrd
665        ['.dwn'] = 1, -- crosswrd
666        ['.ttc'] = 1, -- talk
667        -- ['.txt'] = 1, -- coverpage, but not sure it's safe to include it...
668        ['.eve'] = 1, -- calend0
669        ['.scn'] = 1, -- cwebmac
670        }
671
```

The code to define a specific behaviour for certain files.

```
672
673    luainputenc.unicode_files = {}
674
675    luainputenc.non_unicode_files = {}
676
677    function luainputenc.set_unicode_file(filename)
678        if luainputenc.non_unicode_files[filename] == 1 then
679            luainputenc.non_unicode_files[filename] = nil
680        end
681        luainputenc.unicode_files[filename] = 1
```

```
682     end
683
684     function luainputenc.set_non_unicode_file(filename)
685         if luainputenc.unicode_files[filename] == 1 then
686             luainputenc.unicode_files[filename] = nil
687         end
688         luainputenc.non_unicode_files[filename] = 1
689     end
690
691     function luainputenc.set_unicode_extention(ext)
692         luainputenc.unicode_extention[ext] = 1
693     end
694
695     function luainputenc.set_non_unicode_extention(ext)
696         if luainputenc.unicode_extentions[ext] == 1 then
697             luainputenc.unicode_extentions[ext] = nil
698         end
699     end
700
701     function luainputenc.unset_file(filename)
702         if luainputenc.unicode_files[filename] == 1 then
703             luainputenc.unicode_files[filename] = nil
704         elseif luainputenc.non_unicode_files[filename] == 1 then
705             luainputenc.non_unicode_files[filename] = nil
706         end
707     end
708
709     local unicode, non_unicode = stopped, started
710
711     function luainputenc.find_state(filename)
712         if luainputenc.unicode_files[filename] == 1 then
713             return unicode
714         elseif luainputenc.non_unicode_files[filename] == 1 then
715             return non_unicode
716         else
717             local ext = filename:sub(-4)
718              if luainputenc.unicode_extentions[ext] == 1 then
719                  return unicode
720              else
721                 return non_unicode
722              end
723         end
724     end
725
```

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

```
726
727     function luainputenc.pre_read_file(env)
728         if not env.path then
```

```
729            return
730        end
731        local currentstate = luainputenc.state
732        luainputenc.setstate(luainputenc.find_state(env.filename))
733        env.previousstate = currentstate
734    end
735
736    function luainputenc.close(env)
737        luainputenc.setstate(env.previousstate)
738    end
739
740    luainputenc.callback_registered = 0
741
742    function luainputenc.register_callback()
743        if luainputenc.callback_registered == 0 then
744            callback.add('pre_read_file', luainputenc.pre_read_file,
745                'luainputenc.pre_read_file')
746            callback.add('file_close', luainputenc.close, 'luainputenc.close')
747            luainputenc.callback_registered = 1
748        end
749    end
750
751 end
752
```

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

```
753
754
755 luainputenc.activated_characters = {}
756 luainputenc.characters_are_activated = false
757
758 function luainputenc.declare_character(c)
759    luainputenc.activated_characters[tonumber(c)] = true
760 end
761
762 function luainputenc.force_characters_activated ()
763    luainputenc.characters_are_activated = true
764 end
765
766 function luainputenc.activate_characters()
767    if not luainputenc.characters_are_activated then
768        for n, _ in pairs(luainputenc.activated_characters) do
769            tex.sprint(string.format('\\catcode %d\\active',n))
770        end
771        luainputenc.characters_are_activated = true
772    end
773 end
774
775 function luainputenc.desactivate_characters()
```

```
776      if luainputenc.characters_are_activated then
777          for n, _ in pairs(luainputenc.activated_characters) do
778              tex.sprint(string.format('\\catcode %d=11',n))
779          end
780          luainputenc.characters_are_activated = false
781      end
782 end
783
```