

The `luainputenc` package

Elie Roux

`elie.roux@telecom-bretagne.eu`

2009/11/20 v0.95

Abstract

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

Contents

1 Documentation	1
1.1 Introduction	1
1.2 Overview of 8-bit mode	2
1.3 Overview of UTF-8 mode	2
1.3.1 legacy mode	2
1.3.2 unicode font mode	3
1.3.3 mixed mode	3
2 Files	3
2.1 <code>inputenc.sty</code> patch	3
2.2 <code>luainputenc.sty</code>	4
2.3 <code>lutf8.def</code>	10
2.4 <code>lutf8x.def</code>	12
2.5 <code>eu2enc.def</code>	15
2.6 <code>eu2lmr.fd</code>	15
2.7 <code>luainputenc.lua</code>	16

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 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^AT_EX 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.

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 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 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 `luafontload`.

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 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.

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 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

```

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 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   \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 \LaTeX . 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.

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.

```

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

```

```
141 \repeat{  
142 }
```

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

```
143  
144 \DeclareOption{utf8}{%  
145   \inputencoding{lutf8}%  
146 }  
147  
148 \DeclareOption{lutf8}{%  
149   \inputencoding{lutf8}%  
150 }  
151  
152 \DeclareOption{utf8x}{%  
153   \inputencoding{lutf8}%  
154 }  
155  
156 \DeclareOption{lutf8x}{%  
157   \inputencoding{lutf8x}%  
158 }  
159
```

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

```
160  
161 \DeclareOption{unactivate}{%  
162   \edef\inputencodingname{unactivate}%  
163 }  
164
```

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`.

```
165  
166 \DeclareOption*{  
167   \LIE@activate %  
168   \inpend@loop\^\^\?\^\^ff%  
169   \inputencoding{\CurrentOption}%  
170 }  
171
```

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.

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

```

179 \else
180
    \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.

```

```

181
182 \newcommand*\LIE@setstarted[0]{%
183     \ifnum\LIE@activated=1 %
184         \luadirect{luainputenc.setstarted()}%
185     \fi %
186 }
187
188 \newcommand*\LIE@setstopped[0]{%
189     \ifnum\LIE@activated=1 %
190         \luadirect{luainputenc.setstopped()}%
191     \fi %
192 }
193

```

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 **T_EX**. In **no way** this means you can include true UTF-8 files, it means that you can include files that have been written by **LuaT_EX** 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 **LuaT_EX** and that you want to include.

```

194
195 \newcommand*\LIE@SetUtfFile[1]{%
196     \luadirect{luainputenc.set_unicode_file([[#1]])}%
197 }
198

```

\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).

```

199
200 \newcommand*\LIE@SetNonUtfFile[1]{%
201     \luadirect{luainputenc.set_non_unicode_file([[#1]])}%
202 }
203

```

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

```

204

```

```
205 \newcommand*\lIE@UnsetFile[1]{%
206   \luadirect{luainputenc.unset_file([[#1]])}%
207 }
208
```

\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.

```
209
210 \newcommand*\lIE@SetUtfExt[1]{%
211   \luadirect{luainputenc.set_unicode_extention([[#1]])}%
212 }
213
```

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

```
214
215 \newcommand*\lIE@SetNonUtfExt[1]{%
216   \luadirect{luainputenc.set_non_unicode_extention([[#1]])}%
217 }
218
```

\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.

```
219
220
221 \newcommand*\lIE@InputUtfFile[1]{%
222   \lIE@SetUtfFile{#1}%
223   \input #1%
224   \lIE@UnsetFile{#1}%
225 }
226
```

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

```
227
228 \newcommand*\lIE@InputNonUtfFile[1]{%
229   \lIE@SetNonUtfFile{#1}%
230   \input #1%
231   \lIE@UnsetFile{#1}%
232 }
233
```

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

```
234
235 \newcommand*\InputUtfFile[1]{%
236   \lIE@InputUtfFile{#1}%
237 }
```

```

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

```

2.3 lutf8.def

```

260 %% This file was adapted from utf8.def, which copyright is:
261 %% Copyright 1993 1994 1995 1996 1997 1998 1999 2000 2001 2002 2003
262 %% 2004 2005 2006 The LaTeX3 Project.
263 %%
264 %% utf8.def is under the lppl version 1.3c or later, and can be found
265 %% in the base LaTeX system.
266 %%
267 %% The lppl can be found at http://www.latex-project.org/lppl.txt
268 %%
269 %% The changes to utf8.def are Copyright 2009 Elie Roux, and are under
270 %% the CCO license.
271 %%
272 %% The changes are LuaTeX support.
273 %%
274 %% This file is distributed under the CCO license, with clause 6 of the
275 %% lppl as additional restrictions.
276

```

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.

```

277
278
279 \ProvidesFile{lutf8.def}
280   [2009/11/20 v0.95 UTF-8 support for luainputenc]
281
282 \makeatletter

```

```

283 \catcode`\\ \saved@space@catcode
284
285 \@inpenc@test
286
287 \ifx\@begindocumenthook\@undefined
288   \makeatother
289   \endinput \fi
290

```

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`.

```

291
292 \gdef\DeclareUnicodeCharacter#1#2{%
293   \tempcnta"#1%
294   \catcode\tempcnta\active %
295   \DeclareInputText{\the\tempcnta}{#2}%
296 }
297
298 \onlypreamble\DeclareUnicodeCharacter
299
300 \def\cdp@elt#1#2#3#4{%
301   \wlog{Now handling font encoding #1 ...}%
302   \lowercase{%
303     \InputIfFileExists{#1enc.dfu}}%
304     {\wlog{... processing UTF-8 mapping file for font encoding
305       #1}%
306     \catcode`\_9\relax}%
307     {\wlog{... no UTF-8 mapping file for font encoding #1}}%
308 }
309 \cdp@list
310
311 \def\DeclareFontEncoding@#1#2#3{%
312   \expandafter %
313   \ifx\csname T@#1\endcsname\relax %
314     \def\cdp@elt{\noexpand\cdp@elt}%
315     \xdef\cdp@list{\cdp@list\cdp@elt{#1}%
316       {\default@family}{\default@series}%
317       {\default@shape}}%
318   \expandafter\let\csname#1-cmd\endcsname\@changed@cmd %
319   \begingroup %
320     \wlog{Now handling font encoding #1 ...}%
321     \lowercase{%
322       \InputIfFileExists{#1enc.dfu}}%
323       {\wlog{... processing UTF-8 mapping file for font encoding #1}}%
324       {\wlog{... no UTF-8 mapping file for font encoding #1}}%
325   \endgroup
326   \else
327     \font@info{Redeclaring font encoding #1}%
328   \fi

```

```

329 \global\@namedef{T@\#1}{\#2}%
330 \global\@namedef{M@\#1}{\default@M\#3}%
331 \xdef\LastDeclaredEncoding{\#1}%
332 }
333
334 \DeclareUnicodeCharacter{00A9}{\textcopyright}
335 \DeclareUnicodeCharacter{00AA}{\textordfeminine}
336 \DeclareUnicodeCharacter{00AE}{\textregistered}
337 \DeclareUnicodeCharacter{00BA}{\textordmasculine}
338 \DeclareUnicodeCharacter{02C6}{\textasciicircum}
339 \DeclareUnicodeCharacter{02DC}{\textasciitilde}
340 \DeclareUnicodeCharacter{200C}{\textcompwordmark}
341 \DeclareUnicodeCharacter{2026}{\textellipsis}
342 \DeclareUnicodeCharacter{2122}{\texttrademark}
343 \DeclareUnicodeCharacter{2423}{\textvisiblespace}
344

```

2.4 lutf8x.def

```

345 %% This file was adapted from utf8.def, which copyright is:
346 %% Copyright 1993 1994 1995 1996 1997 1998 1999 2000 2001 2002 2003
347 %% 2004 2005 2006 The LaTeX3 Project.
348 %%
349 %% utf8.def is under the lppl version 1.3c or later, and can be found
350 %% in the base LaTeX system.
351 %%
352 %% The lppl can be found at http://www.latex-project.org/lppl.txt
353 %%
354 %% The changes to utf8.def are Copyright 2009 Elie Roux, and are under
355 %% the CCO license.
356 %%
357 %% The changes are LuaTeX support.
358 %%
359 %% This file is distributed under the CCO license, with clause 6 of the
360 %% lppl as additional restrictions.
361

```

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`.

```

362
363 \ProvidesFile{lutf8x.def}
364   [2009/11/20 v0.95 UTF-8 support for luainputenc]
365
366 \makeatletter
367 \catcode`\ \saved@space@catcode
368
369 \cinpenc@test
370
371 \ifx\@begindocumenthook\undefined

```

```

372 \makeatother
373 \endinput \fi
374

```

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

```

375
376 \gdef\DeclareUnicodeCharacter#1#2{%
377   \tempcnta"#1%
378   \luadirect{luainputenc.declare_character('the\tempcnta')}%
379   \catcode@\tempcnta\active %
380   \DeclareInputText{\the\tempcnta}{#2}%
381 }
382
383 \onlypreamble\DeclareUnicodeCharacter
384
385 \def\cdp@elt#1#2#3#4{%
386   \wlog{Now handling font encoding #1 ...}%
387   \lowercase{%
388     \InputIfFileExists{#1enc.dfu}}%
389     {\wlog{... processing UTF-8 mapping file for font encoding
390       #1}%
391     \catcode`\_9\relax}%
392     {\wlog{... no UTF-8 mapping file for font encoding #1}}%
393 }
394 \cdp@list
395

```

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

```

396
397 \def\lIE@ActivateUnicodeCatcodes{%
398 \luadirect{luainputenc.activate_characters()}%
399 }
400
401 \def\lIE@DesactivateUnicodeCatcodes{%
402 \luadirect{luainputenc.desactivate_characters()}%
403 }
404
405 \def\lIE@CharactersActivated{%
406 \luadirect{luainputenc.force_characters_activated()}%
407 }
408
409 \edef\lIE@EU{EU2}
410

```

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`.

```

411
412 \def\DeclareFontEncoding@#1#2#3{%
413   \edef\lIE@test{#1}%
414   \ifx\lIE@test\lIE@EU %

```

```

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

```

```

465 \DeclareUnicodeCharacter{2026}{\textellipsis}
466 \DeclareUnicodeCharacter{2122}{\texttrademark}
467 \DeclareUnicodeCharacter{2423}{\textvisiblespace}
468

```

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 L^AT_EX 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 `luatofont` is a good choice to be able to do so.

```

469
470 \ProvidesFile{eu2enc.def}[2009/11/20 v0.1 a unicode font encoding for LuaTeX.]
471 \DeclareFontEncoding{EU2}{}{}
472 \DeclareErrorFont{EU2}{lmr}{m}{n}{10}
473 \DeclareFontSubstitution{EU2}{lmr}{m}{n}
474

```

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.

```

475
476 \ProvidesFile{eu2lmr.fd}
477     [2009/11/20 v0.2 Font defs for Latin Modern for LuaTeX's EU2 encoding]
478 \DeclareFontFamily{EU2}{lmr}{}%
479 \DeclareFontShape{EU2}{lmr}{m}{n}%
480     {<-5.5> "lmroman5-regular:+tlig;+tsub;+liga;+rlig;"}
481     <5.5-6.5> "lmroman6-regular:+tlig;+tsub;+liga;+rlig;"}
482     <6.5-7.5> "lmroman7-regular:+tlig;+tsub;+liga;+rlig;"}
483     <7.5-8.5> "lmroman8-regular:+tlig;+tsub;+liga;+rlig;"}
484     <8.5-9.5> "lmroman9-regular:+tlig;+tsub;+liga;+rlig;"}
485     <9.5-11> "lmroman10-regular:+tlig;+tsub;+liga;+rlig;"}
486     <11-15> "lmroman12-regular:+tlig;+tsub;+liga;+rlig;"}
487     <15-> "lmroman17-regular:+tlig;+tsub;+liga;+rlig;"}
488     }{}%
489 \DeclareFontShape{EU2}{lmr}{m}{sl}%
490     {<-8.5> "lmroman8-oblique:+tlig;+tsub;+liga;+rlig;"}
491     <8.5-9.5> "lmroman9-oblique:+tlig;+tsub;+liga;+rlig;"}
492     <9.5-11> "lmroman10-oblique:+tlig;+tsub;+liga;+rlig;"}
493     <11-15> "lmroman12-oblique:+tlig;+tsub;+liga;+rlig;"}
494     <15-> "lmroman17-oblique:+tlig;+tsub;+liga;+rlig;"}
495     }{}%
496 \DeclareFontShape{EU2}{lmr}{m}{it}%
497     {<-7.5> "lmroman7-italic:+tlig;+tsub;+liga;+rlig;"}
498     <7.5-8.5> "lmroman8-italic:+tlig;+tsub;+liga;+rlig;"}
499     <8.5-9.5> "lmroman9-italic:+tlig;+tsub;+liga;+rlig;"}
500     <9.5-11> "lmroman10-italic:+tlig;+tsub;+liga;+rlig;"}
501     <11-> "lmroman12-italic:+tlig;+tsub;+liga;+rlig;"}
502     }{}%

```

```

503 \DeclareFontShape{EU2}{lmr}{m}{sc}%
504     {<-> "lmroman10-capsregular:+tlig;+tsub;+liga;+rlig;"{}}
505 %
506 % Is this the right 'shape'?:
507 \DeclareFontShape{EU2}{lmr}{m}{scsl}%
508     {<-> "lmroman10-capsoblique:+tlig;+tsub;+liga;+rlig;"{}}
509 %%%%%% bold series
510 \DeclareFontShape{EU2}{lmr}{b}{n}%
511     {<-> "lmroman10-demi:+tlig;+tsub;+liga;+rlig;"{}}
512 \DeclareFontShape{EU2}{lmr}{b}{sl}%
513     {<-> "lmroman10-demioblique:+tlig;+tsub;+liga;+rlig;"{}}
514 %%%%%% bold extended series
515 \DeclareFontShape{EU2}{lmr}{bx}{n}%
516     {<-5.5> "lmroman5-bold:+tlig;+tsub;+liga;+rlig;"%
517     <5.5-6.5> "lmroman6-bold:+tlig;+tsub;+liga;+rlig;"%
518     <6.5-7.5> "lmroman7-bold:+tlig;+tsub;+liga;+rlig;"%
519     <7.5-8.5> "lmroman8-bold:+tlig;+tsub;+liga;+rlig;"%
520     <8.5-9.5> "lmroman9-bold:+tlig;+tsub;+liga;+rlig;"%
521     <9.5-11> "lmroman10-bold:+tlig;+tsub;+liga;+rlig;"%
522     <11-> "lmroman12-bold:+tlig;+tsub;+liga;+rlig;"%
523     }{}}
524 \DeclareFontShape{EU2}{lmr}{bx}{it}%
525     {<-> "lmroman10-bolditalic:+tlig;+tsub;+liga;+rlig;"{}}
526 \DeclareFontShape{EU2}{lmr}{bx}{sl}%
527     {<-> "lmroman10-boldoblique:+tlig;+tsub;+liga;+rlig;"{}}
528

```

2.7 luainputenc.lua

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

```

529
530 luainputenc = { }
531
532 luainputenc.module = {
533     name      = "luainputenc",
534     version   = 0.95,
535     date      = "2009/11/20",
536     description = "Lua simple inputenc package.",
537     author    = "Elie Roux",
538     copyright = "Elie Roux",
539     license   = "CC0",
540 }
541
542 luatextra.provides_module(luainputenc.module)
543
544 local format = string.format
545
546 luainputenc.log = luainputenc.log or function(...)
547     luatextra.module_log('luainputenc', format(...))
548 end

```

```

549
550 local char, utfchar, byte, format, gsub, utfbyte, utfgsub =
551 string.char, unicode.utf8.char, string.byte, string.format, string.gsub, unicode.utf8.byte, unicode.utf8.gsub
552

```

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

```

553
554 function luainputenc.byte_to_utf(ch)
555     return utfchar(byte(ch))
556 end
557

```

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

```

558
559 function luainputenc.fake_utf_read(buf)
560     return gsub(buf,"(.)", luainputenc.byte_to_utf)
561 end
562

```

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

```

563
564 function luainputenc.utf_to_byte(ch)
565     return char(utfbyte(ch))
566 end
567

```

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

```

568
569 function luainputenc.fake_utf_write(buf)
570     return utfgsub(buf,"(.)", luainputenc.utf_to_byte)
571 end
572

```

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.

```

573
574 if tex.luatexversion > 42 then
575
576     function luainputenc.register_callbacks()
577         callback.add('process_output_buffer', luainputenc.fake_utf_write, 'luainputenc.fake_utf_write')
578         callback.add('process_input_buffer', luainputenc.fake_utf_read, 'luainputenc.fake_utf_read')
579     end
580
581 else
582

```

`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.

```

583     local started, stopped = 1, 0
585
586     luainputenc.state = stopped
587
588     function luainputenc.setstate(state)
589         if state == luainputenc.state then
590             return
591         elseif state == started then
592             luainputenc.start()
593         else
594             luainputenc.stop()
595         end
596     end
597
598     function luainputenc.setstarted()
599         luainputenc.setstate(started)
600     end
601
602     function luainputenc.setstopped()
603         luainputenc.setstate(stopped)
604     end
605
606     function luainputenc.start()
607         callback.add('process_input_buffer', luainputenc.fake_utf_read,
608                     'luainputenc.fake_utf_read')
609         luainputenc.state = started
610         if luainputenc.callback_registered == 0 then
611             luainputenc.register_callback()
612         end
613     end
614
615     function luainputenc.stop()
616         callback.remove('process_input_buffer', 'luainputenc.fake_utf_read')
617         luainputenc.state = stopped
618         return
619     end
620

```

Here is a list of all file extenstions for which we consider that the files have been written by LuaT_EX, 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.

```

621
622     luainputenc.unicode_extentions = {
623         ['.aux'] = 1, -- basic files
624         ['.toc'] = 1,
625         ['.gls'] = 1,
626         ['.ind'] = 1,
627         ['.idx'] = 1,
628         ['.vrb'] = 1, -- beamer and powerdot

```

```

629      ['.nav'] = 1, -- other beamer extention
630      ['.sol'] = 1,
631      ['.qsl'] = 1,
632      ['.snm'] = 1,
633      ['.pgn'] = 1, -- pagereference
634      ['.cpg'] = 1, -- AlProTeX
635      ['.pst'] = 1, -- pst-tree
636      ['.tmp'] = 1, -- sauerj/collect
637      ['.sym'] = 1, -- listofsymbols
638      ['.sub'] = 1, -- listofsymbols
639      ['.lof'] = 1, -- preprint
640      ['.lot'] = 1, -- preprint
641      ['.mtc1'] = 1, -- minitoc
642      ['.ovr'] = 1, -- thumbss
643      ['.fff'] = 1, -- endplate
644      ['.sbb'] = 1, -- splitbib
645      ['.bb1'] = 1, -- latex
646      ['.ain'] = 1, -- authorindex
647      ['.abb'] = 1, -- juraabbrev
648      ['.ent'] = 1, -- endnotes
649      ['.end'] = 1, -- fn2end
650      ['.thm'] = 1, -- ntheorem
651      ['.xtr'] = 1, -- extract
652      ['.han'] = 1, -- lingoaho
653      ['.bnd'] = 1, -- bibref
654      ['.bb1'] = 1, -- bibref
655      ['.col'] = 1, -- mwrite
656      ['.ttt'] = 1, -- endfloat
657      ['.fax'] = 1, -- lettre
658      ['.tns'] = 1, -- lettre
659      ['.odt'] = 1, -- lettre
660      ['.etq'] = 1, -- lettre
661      ['.emd'] = 1, -- poemscol
662      ['.emx'] = 1, -- poemscol
663      ['.ctn'] = 1, -- poemscol
664      ['.hst'] = 1, -- vhstory
665      ['.acr'] = 1, -- crosswrd
666      ['.dwn'] = 1, -- crosswrd
667      ['.ttc'] = 1, -- talk
668      -- ['.txt'] = 1, -- coverpage, but not sure it's safe to include it...
669      ['.eve'] = 1, -- calend0
670      ['.scn'] = 1, -- cwebmac
671      }
672

```

The code to define a specific behaviour for certain files.

```

673
674      luainputenc_unicode_files = {}
675
676      luainputenc_non_unicode_files = {}
677

```

```

678     function luainputenc.set_unicode_file(filename)
679         if luainputenc.non_unicode_files[filename] == 1 then
680             luainputenc.non_unicode_files[filename] = nil
681         end
682         luainputenc_unicode_files[filename] = 1
683     end
684
685     function luainputenc.set_non_unicode_file(filename)
686         if luainputenc_unicode_files[filename] == 1 then
687             luainputenc_unicode_files[filename] = nil
688         end
689         luainputenc.non_unicode_files[filename] = 1
690     end
691
692     function luainputenc.set_unicode_extention(ext)
693         luainputenc_unicode_extention[ext] = 1
694     end
695
696     function luainputenc.set_non_unicode_extention(ext)
697         if luainputenc_unicode_extentions[ext] == 1 then
698             luainputenc_unicode_extentions[ext] = nil
699         end
700     end
701
702     function luainputenc.unset_file(filename)
703         if luainputenc_unicode_files[filename] == 1 then
704             luainputenc_unicode_files[filename] = nil
705         elseif luainputenc_non_unicode_files[filename] == 1 then
706             luainputenc_non_unicode_files[filename] = nil
707         end
708     end
709
710     local unicode, non_unicode = stopped, started
711
712     function luainputenc.find_state(filename)
713         if luainputenc_unicode_files[filename] == 1 then
714             return unicode
715         elseif luainputenc_non_unicode_files[filename] == 1 then
716             return non_unicode
717         else
718             local ext = filename:sub(-4)
719             if luainputenc_unicode_extentions[ext] == 1 then
720                 return unicode
721             else
722                 return non_unicode
723             end
724         end
725     end
726

```

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

```
727
728     function luainputenc.pre_read_file(env)
729         if not env.path then
730             return
731         end
732         local currentstate = luainputenc.state
733         luainputenc.setstate(luainputenc.find_state(env.filename))
734         env.previousstate = currentstate
735     end
736
737     function luainputenc.close(env)
738         luainputenc.setstate(env.previousstate)
739     end
740
741     luainputenc.callback_registered = 0
742
743     function luainputenc.register_callback()
744         if luainputenc.callback_registered == 0 then
745             callback.add('pre_read_file', luainputenc.pre_read_file,
746                         'luainputenc.pre_read_file')
747             callback.add('file_close', luainputenc.close, 'luainputenc.close')
748             luainputenc.callback_registered = 1
749         end
750     end
751
752 end
753
```

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

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

```
772         luainputenc.characters_are_activated = true
773     end
774 end
775
776 function luainputenc.desactivate_characters()
777     if luainputenc.characters_are_activated then
778         for n, _ in pairs(luainputenc.activated_characters) do
779             tex.sprint(string.format('\\catcode %d=11',n))
780         end
781         luainputenc.characters_are_activated = false
782     end
783 end
784
```