

The `hyphen.cfg` file for Lua_T_EX

Khaled Hosny, Élie Roux, and Manuel Pégourié-Gonnard

`khaledhosny@eglug.org`

`elie.roux@telecom-bretagne.eu`

`mpg@elzevir.fr`

2010/04/28 v1.3beta

Abstract

This is a modified version of the file `hyphen.cfg` distributed with the `babel` package, with a supporting Lua module, aimed at adapting `babel`'s hyphenation patterns loading mechanism to Lua_T_EX's dynamic pattern loading capabilities. It makes use of a `language.dat.lua` file (whose format is described below) that should be present in the distribution, in addition to the regular `language.dat` file.

Much of the modified code here is shared with a version of `etex.src` modified for the same reasons, which also make use of the `luatex-hyphen.lua` file described here.

1 Documentation

Hyphenation patterns should be loaded at runtime with Lua_T_EX: if they appear in the format, they will be rehashed when the format is loaded anyway, which makes the format quite long to load (many seconds even on modern machines) and provides for bad user experience. Hence, it is desirable to load as few patterns as possible in the format, and load on-demand the needed patterns at runtime.

This package provides a modified version of `hyphen.cfg` adapted to Lua_T_EX, as well as a supporting Lua module. Since a lot of things, especially the catcodes, are not as predictable at runtime than at format creation time, we don't `\input` the usual pattern files, but rather load the patterns using the Lua interface, using a special plain text version of the pattern files if available.

The existence and file name of such a version cannot be guessed, so we need a specific database: the file `language.dat.lua`. This file should be loadable by Lua and return a table whose keys are the canonical language names as found in `language.dat`, and the values are Lua tables consisting of:

1. A fixed part with two fields:

```
loader = <string> name of the TeX loader
synonyms = { <string> alternative name
```

The `loader` field is currently unused.

2. A variable part consisting of either:

- For most languages:

```
code = <string> language code
lefthyphenmin = <number> value for \letfhyphenmin
righthyphenmin = <number> value for \letfhyphenmin
```

The `code` field determines where patterns and exceptions will be searched: in files `hyph-<code>.pat.txt` and `hyph-<code>.hyp.txt` respectively. The values of `*hyphenmin` are currently unused.

- Special case are supported by a field `special`. Currently, two kinds of value are recognized: `'null'` for languages with no hyphenation patterns nor exceptions. `'disabled:<reason>'` allows to disable specific languages: when the user tries to load this language, an error will be issued, with the `<reason>`.

Languages that are mentioned in `language.dat` but not in `language.dat.lua` will be loaded in the format. So, if the `language.dat.lua` file is missing or incomplete, languages will just go back to the “old” behaviour, resulting in longer startup time, which does seem less bad than complete breakage.

For backward compatibility, Knuth’s original patterns for US English are always loaded in the format, as `\language0`.

The modified version of `hyphen.cfg` provided here checks for the engine, and should continue to work with any engine without any modified behaviour. However, it is recommended to install it in such a way that the original `hyphen.cfg` from `babel` is found first by any engine other than `LuaTeX`.

2 Package code

2.1 luatex-hyphen.lua

```
1
2 luatexhyphen = {}
3
4 luatexhyphen.version = "1.3beta"
5
6 local dbname = "language.dat.lua"
7
8 local function warn (msg, ...)
9     texio.write_nl('luatex-hyphen: '..string.format(msg, ...))
10 end
11
12 luatexhyphen.language_dat = {}
13 local dbfile = kpse.find_file(dbname)
14 if not dbfile then
15     warn("file not found: "..dbname)
```

```

16 else
17     luatexhyphen.language_dat = dofile(dbfile)
18 end
19
20 function luatexhyphen.lookupname(l)
21     if luatexhyphen.language_dat[l] then
22         return luatexhyphen.language_dat[l], l
23     else
24         for orig,lt in pairs(luatexhyphen.language_dat) do
25             for _,syn in ipairs(lt.synonyms) do
26                 if syn == l then
27                     return lt, orig
28                 end
29             end
30         end
31     end
32     return nil
33 end
34
35 function luatexhyphen.loadlanguage(l, id)
36     local lt, orig = luatexhyphen.lookupname(l)
37     if not lt then
38         warn("no entry in %s for this language: %s", dbname, l)
39         return
40     end
41     local msg = "loading%s patterns and exceptions for: %s (\\language%d)"
42     if lt.special then
43         if lt.special == 'null' then
44             warn(msg, ' (null)', orig, id)
45         elseif lt.special:find('^disabled:') then
46             warn("language disabled by %s: %s (%s)", dbname, orig,
47                 lt.special:gsub('^disabled:', ''))
48         else
49             warn("bad entry in %s for language %s")
50         end
51         return
52     end
53     warn(msg, '', orig, id)
54     for ext, fun in pairs({pat = lang.patterns, hyp = lang.hyphenation}) do
55         local n = 'hyph-'..lt.code..'..'..ext..'..txt'
56         local f = kpse.find_file(n)
57         if not f then
58             warn("file not found: %s", n)
59             return
60         end
61         f = io.open(f, 'r')
62         local data = f:read('*a')
63         f:close()
64         if not data then
65             warn("file not readable: %s", f)

```

```

66         return
67     end
68     fun(lang.new(id), data)
69 end
70 end

```

2.2 hyphen.cfg

```

71 \ifx\ProvidesFile\@undefined
72 \def\ProvidesFile#1[#2 #3 #4]{%
73 \wlog{File: #1 #4 #3 <#2>}%

    Use a modified banner for LuaTeX.
74 \ifx\directlua\@undefined
75 \toks8{Babel <#3> and hyphenation patterns for }%
76 \else
77 \toks8{LuaTeX adaptation of babel <#3>
78 and hyphenation patterns for }%
79 \fi

80 \let\ProvidesFile\@undefined
81 }
82 \def\ProvidesLanguage#1[#2 #3 #4]{%
83 \wlog{Language: #1 #4 #3 <#2>}%
84 }
85 \else
86 \let\bbl@tempa\ProvidesFile
87 \def\ProvidesFile#1[#2 #3 #4]{%

    Same here.
88 \ifx\directlua\@undefined
89 \toks8{Babel <#3> and hyphenation patterns for }%
90 \else
91 \toks8{LuaTeX adaptation of babel <#3>
92 and hyphenation patterns for }%
93 \fi

94 \bbl@tempa#1[#2 #3 #4]%
95 \let\ProvidesFile\bbl@tempa}
96 \def\ProvidesLanguage#1{%
97 \begingroup
98 \catcode'\ 10 %
99 \@makeother\/%
100 \@ifnextchar[%]
101 {\@provideslanguage{#1}}{\@provideslanguage{#1}[]}}
102 \def\@provideslanguage#1[#2]{%
103 \wlog{Language: #1 #2}%
104 \expandafter\xdef\csname ver@#1.1df\endcsname{#2}%
105 \endgroup}
106 \fi
107

```

File identification is modified again.

```

108 \ProvidesFile{hyphen.cfg}
109         [2010/04/26 v3.81-luatex-1.3beta %
110         Language switching mechanism for LuaTeX, adapted from babel v3.81]

111 \ifx\AtBeginDocument\@undefined
112 \input plain.def\relax
113 \fi
114 \ifx\language\@undefined
115 \csname newcount\endcsname\language
116 \fi
117 \ifx\newlanguage\@undefined
118 \csname newcount\endcsname\last@language
119 \else
120 \countdef\last@language=19
121 \fi
122 \ifx\newlanguage\@undefined
123 \def\addlanguage#1{%
124 \global\advance\last@language \@ne
125 \ifnum\last@language<\@cclvi
126 \else
127 \errmessage{No room for a new \string\language!}%
128 \fi
129 \global\chardef#1\last@language
130 \wlog{\string#1 = \string\language\the\last@language}}
131 \else
132 \def\addlanguage{\alloc@9\language\chardef\@cclvi}
133 \fi
134 \def\adddialect#1#2{%
135 \global\chardef#1#2\relax
136 \wlog{\string#1 = a dialect from \string\language#2}}
137 \def\iflanguage#1{%
138 \expandafter\ifx\csname l@#1\endcsname\relax
139 \@nolanerr{#1}%
140 \else
141 \bbl@afterfi{\ifnum\csname l@#1\endcsname=\language
142 \expandafter\@firstoftwo
143 \else
144 \expandafter\@secondoftwo
145 \fi}%
146 \fi}
147 \edef\selectlanguage{%
148 \noexpand\protect
149 \expandafter\noexpand\csname selectlanguage \endcsname
150 }
151 \ifx\@undefined\protect\let\protect\relax\fi
152 \ifx\documentclass\@undefined
153 \def\xstring{\string\string\string}
154 \else
155 \let\xstring\string
156 \fi

```

```

157 \xdef\bbl@language@stack{}
158 \def\bbl@push@language{%
159   \xdef\bbl@language@stack{\language+\bbl@language@stack}%
160 }
161 \def\bbl@pop@lang#1+#2-#3{%
162   \def\language{#1}\xdef#3{#2}%
163 }
164 \def\bbl@pop@language{%
165   \expandafter\bbl@pop@lang\bbl@language@stack-\bbl@language@stack
166   \expandafter\bbl@set@language\expandafter{\language}%
167 }
168 \expandafter\def\csname selectlanguage \endcsname#1{%
169   \bbl@push@language
170   \aftergroup\bbl@pop@language
171   \bbl@set@language{#1}}
172 \def\bbl@set@language#1{%
173   \edef\language{%
174     \ifnum\escapechar=\expandafter'\string#1\@empty
175     \else \string#1\@empty\fi}%
176   \select@language{\language}%
177   \if@files
178     \protected@write\@auxout{\string\select@language{\language}}%
179     \addtocontents{toc}{\xstring\select@language{\language}}%
180     \addtocontents{lof}{\xstring\select@language{\language}}%
181     \addtocontents{lot}{\xstring\select@language{\language}}%
182   \fi}
183 \def\select@language#1{%
184   \expandafter\ifx\csname l@#1\endcsname\relax
185     \@nolanerr{#1}%
186   \else
187     \expandafter\ifx\csname date#1\endcsname\relax
188       \@noopterr{#1}%
189     \else
190       \bbl@patterns{\language}%
191       \originalTeX
192       \expandafter\def\expandafter\originalTeX
193         \expandafter{\csname noextras#1\endcsname
194           \let\originalTeX\@empty}%
195       \languageshorthands{none}%
196       \babel@beginsave
197       \csname captions#1\endcsname
198       \csname date#1\endcsname
199       \csname extras#1\endcsname\relax
200       \babel@savevariable\lefthyphenmin
201       \babel@savevariable\righthyphenmin
202       \expandafter\ifx\csname #1hyphenmins\endcsname\relax
203         \set@hyphenmins\tw@\thr@\relax
204       \else
205         \expandafter\expandafter\expandafter\set@hyphenmins
206         \csname #1hyphenmins\endcsname\relax

```

```

207     \fi
208     \fi
209 \fi}
210 \long\def\otherlanguage#1{%
211   \csname selectlanguage \endcsname{#1}%
212   \ignorespaces
213 }
214 \long\def\endotherlanguage{%
215   \originalTeX
216   \global\@ignoretrue\ignorespaces
217 }
218 \expandafter\def\csname otherlanguage*\endcsname#1{%
219   \foreign@language{#1}%
220 }
221 \expandafter\def\csname endotherlanguage*\endcsname{%
222   \csname noextras\language\endcsname
223 }
224 \def\foreignlanguage{\protect\csname foreignlanguage \endcsname}
225 \expandafter\def\csname foreignlanguage \endcsname#1#2{%
226   \begingroup
227     \originalTeX
228     \foreign@language{#1}%
229     #2%
230     \csname noextras#1\endcsname
231   \endgroup
232 }
233 \def\foreign@language#1{%
234   \def\language{#1}%
235   \expandafter\ifx\csname l@#1\endcsname\relax
236     \nol@nerr{#1}%
237   \else
238     \bbl@patterns{\language}%
239     \languageshortands{none}%
240     \csname extras#1\endcsname
241     \expandafter\ifx\csname #1hyphenmins\endcsname\relax
242       \set@hyphenmins\tw@\thr@@\relax
243     \else
244       \expandafter\expandafter\expandafter\set@hyphenmins
245       \csname #1hyphenmins\endcsname\relax
246     \fi
247   \fi
248 }
249 \def\bbl@patterns#1{%
250   \language=\expandafter\ifx\csname l@#1:\f@encoding\endcsname\relax
251     \csname l@#1\endcsname
252   \else
253     \csname l@#1:\f@encoding\endcsname
254   \fi\relax

```

With LuaTeX, load patterns and exceptions on the fly using functions from the

supporting Lua module, unless of course they are already loaded for this language (identified by its number to avoid problems with synonyms).

Also, since this code will be executed at runtime, be careful while testing if we're using LuaTeX.

```

255 \ifx\directlua\@undefined\else
256 \ifx\directlua\relax\else
257 \ifcsname lu@texhyphen@loaded@\the\language\endcsname \else
258 \global\@namedef{lu@texhyphen@loaded@\the\language}{}%
259 \directlua{
260     if not luatexhyphen then
261         dofile(assert(kpse.find_file("luatex-hyphen.lua")))
262     end
263     luatexhyphen.loadlanguage("\luatexluaescapestring{#1}",
264         \the\language)}%
265 \fi
266 \fi
267 \fi

268 }
269 \def\hyphenrules#1{%
270 \expandafter\ifx\csname l@#1\endcsname\@undefined
271 \nolanerr{#1}%
272 \else
273 \bbl@patterns{#1}%
274 \languageshorthands{none}%
275 \expandafter\ifx\csname #1hyphenmins\endcsname\relax
276 \set@hyphenmins\tw@\thr@@\relax
277 \else
278 \expandafter\expandafter\expandafter\set@hyphenmins
279 \csname #1hyphenmins\endcsname\relax
280 \fi
281 \fi
282 }
283 \def\endhyphenrules{}
284 \def\providehyphenmins#1#2{%
285 \expandafter\ifx\csname #1hyphenmins\endcsname\relax
286 \@namedef{#1hyphenmins}{#2}%
287 \fi}
288 \def\set@hyphenmins#1#2{\lefthyphenmin#1\righthyphenmin#2}
289 \def\LdfInit{%
290 \chardef\atcatcode=\catcode'\@
291 \catcode'\@=11\relax
292 \input babel.def\relax
293 \catcode'\@=\atcatcode \let\atcatcode\relax
294 \LdfInit}
295 \ifx\originalTeX\@undefined\let\originalTeX\@empty\fi
296 \ifx\babel@beginsave\@undefined\let\babel@beginsave\relax\fi
297 \ifx\PackageError\@undefined
298 \def\@nolanerr#1{%
299 \errhelp{Your command will be ignored, type <return> to proceed}%

```



```

300   \errmessage{You haven't defined the language #1\space yet}}
301 \def\@nopatterns#1{%
302   \message{No hyphenation patterns were loaded for}%
303   \message{the language '#1'}%
304   \message{I will use the patterns loaded for \string\language=0
305     instead}}
306 \def\@noopterr#1{%
307   \errmessage{The option #1 was not specified in \string\usepackage}
308   \errhelp{You may continue, but expect unexpected results}}
309 \def\@activated#1{%
310   \wlog{Package babel Info: Making #1 an active character}}
311 \else
312   \newcommand*\@nolanerr}[1]{%
313     \PackageError{babel}%
314       {You haven't defined the language #1\space yet}%
315       {Your command will be ignored, type <return> to proceed}}
316   \newcommand*\@nopatterns}[1]{%
317     \PackageWarningNoLine{babel}%
318       {No hyphenation patterns were loaded for\MessageBreak
319         the language '#1'\MessageBreak
320         I will use the patterns loaded for \string\language=0
321         instead}}
322   \newcommand*\@noopterr}[1]{%
323     \PackageError{babel}%
324       {You haven't loaded the option #1\space yet}%
325       {You may proceed, but expect unexpected results}}
326   \newcommand*\@activated}[1]{%
327     \PackageInfo{babel}{%
328       Making #1 an active character}}
329 \fi
330 \def\process@line#1#2 #3/{%
331   \ifx=#1
332     \process@synonym#2 /
333   \else
334     \process@language#1#2 #3/%
335   \fi
336 }
337 \toks@{}
338 \def\process@synonym#1 /{%
339   \ifnum\last@language=\m@ne
340     \expandafter\chardef\csname l@#1\endcsname\relax
341     \wlog{\string\l@#1=\string\language0}
342     \toks@\expandafter\the\toks@
343     \expandafter\let\csname #1hyphenmins\expandafter\endcsname
344     \csname\language\hyphenmins\endcsname}%
345   \else
346     \expandafter\chardef\csname l@#1\endcsname\last@language
347     \wlog{\string\l@#1=\string\language\the\last@language}
348     \expandafter\let\csname #1hyphenmins\expandafter\endcsname
349     \csname\language\hyphenmins\endcsname

```

```

350 \fi
351 }
352 \def\process@language#1 #2 #3/{%
353 \expandafter\addlanguage\csname l@#1\endcsname
354 \expandafter\language\csname l@#1\endcsname
355 \def\language#1}%

In the LuaTeX case, we have to decide whether to load the language now.
356 \ifx\directlua\@undefined
357 \global\toks8\expandafter{\the\toks8#1, }%
358 \else
359 \directlua{
360 if not luatexhyphen then
361 dofile(assert(kpse.find_file("luatex-hyphen.lua")))
362 end
363 processnow = (tex.language == 0) or
364 (luatexhyphen.lookupname("\luatexluaescapestring{#1}") == nil)}%
365 \ifnum0=\directlua{tex.sprint(processnow and "0" or "1")}\relax
366 \global\toks8\expandafter{\the\toks8#1, }%
367 \global\@namedef{lu@texhyphen@loaded@the\language}{}%
368 \fi
369 \fi

370 \begingroup
371 \bbl@get@enc#1:@@@
372 \ifx\bbl@hyph@enc@empty
373 \else
374 \fontencoding{\bbl@hyph@enc}\selectfont
375 \fi
376 \lefthyphenmin\m@ne

```

Assume the first (that is, zeroth) language in `language.dat` is English. This assumption is very reasonable, since otherwise it would break compatibility with frozen TeX by not providing Knuth's original patterns as `\language0`, so we're pretty sure about this point. We do load this first language, since we want Knuth's patterns to be active as soon as the format is loaded.

```

377 \ifx\directlua\@undefined
378 \input #2\relax
379 \else
380 \ifnum0=\directlua{tex.sprint(processnow and "0" or "1")}\relax
381 \input #2\relax
382 \fi
383 \fi

384 \ifnum\lefthyphenmin=\m@ne
385 \else
386 \expandafter\xdef\csname #1hyphenmins\endcsname{%
387 \the\lefthyphenmin\the\righthyphenmin}%
388 \fi
389 \endgroup
390 \ifnum\the\language=\z@
391 \expandafter\ifx\csname #1hyphenmins\endcsname\relax

```

```

392     \set@hyphenmins\tw@\thr@@\relax
393   \else
394     \expandafter\expandafter\expandafter\set@hyphenmins
395       \csname #1hyphenmins\endcsname
396   \fi
397   \the\toks@
398 \fi
399 \toks@{}%
400 \def\bbl@tempa{#3}%
401 \ifx\bbl@tempa\@empty
402 \else
403   \ifx\bbl@tempa\space
404   \else

```

Likewise, don't load hyphenation exceptions now, but rather when we load the patterns. (Anyway, in practice, the third field of `language.dat` is never used since exceptions are defined in the same file as patterns, so it doesn't really matter.)

There are no hyphenation exceptions for english, and since it is frozen, we can rely on this, so no need for a special case for english here.

```

405   \ifx\directlua\@undefined
406     \input #3\relax
407   \else
408     \ifnum0=\directlua{tex.sprint(processnow and "0" or "1")}\relax
409       \input #3\relax
410     \fi
411     \directlua{processnow = nil}%
412   \fi
413 \fi
414 \fi
415 }
416 \def\bbl@get@enc#1:#2\@@@{%
417   \def\bbl@tempa{#1}%
418   \def\bbl@tempb{#2}%
419   \ifx\bbl@tempb\@empty
420     \let\bbl@hyph@enc\@empty
421   \else
422     \bbl@get@enc#2\@@@
423     \edef\bbl@hyph@enc{\bbl@tempa}%
424   \fi}
425 \openin1 = language.dat
426 \ifeof1
427   \message{I couldn't find the file language.dat,\space
428     I will try the file hyphen.tex}
429   \input hyphen.tex\relax
430 \else
431   \last@language\m@ne
432   \loop
433     \endlinechar\m@ne
434     \read1 to \bbl@line
435     \endlinechar'\^M

```

```

436 \ifx\bbl@line\@empty
437 \else
438 \edef\bbl@line{\bbl@line\space/}%
439 \expandafter\process@line\bbl@line
440 \fi
441 \iftrue \csname fi\endcsname
442 \csname if\ifeof1 false\else true\fi\endcsname
443 \repeat
444 \language=0
445 \fi
446 \closein1
447 \let\process@language\@undefined
448 \let\process@synonym\@undefined
449 \let\process@line\@undefined
450 \let\bbl@tempa\@undefined
451 \let\bbl@tempb\@undefined
452 \let\bbl@eq\@undefined
453 \let\bbl@line\@undefined
454 \let\bbl@get@enc\@undefined
455 \ifx\addto@hook\@undefined
456 \else
457 \expandafter\addto@hook\expandafter\everyjob\expandafter{%
458 \expandafter\typeout\expandafter{\the\toks8 loaded.}}
459 \fi

```