

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, ... }
```

Those field's values must be the same as in `language.dat`. 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, the following kind of values 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>`.

`0` only `english` should use this type of special, to indicate it is normally dumped in the format (see below).

Special languages may have `*hyphenmin` information when it makes sense (mostly `\language0`).

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 seems less bad than complete breakage.

For backward compatibility, Knuth's original patterns for US English are always loaded in the format, as `\language0`.¹

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 <lua>
```

Start a Lua module, importing only the necessary functions as locals.

```
2 local error, dofile, pairs, ipairs = error, dofile, pairs, ipairs
3 local io, texio, lang, kpse = io, texio, lang, kpse
4 module('luatexhyphen')
```

¹It is assumed to be the first entry in `language.dat`.

Two functions for error and information reporting.

```

5 local function wlog(msg, ...)
6     texio.write_nl('log', 'luatex-hyphen: '..msg:format(...))
7 end
8 local function err(msg, ...)
9     error('luatex-hyphen: '..msg:format(...), 2)
10 end

```

Load the language.dat.lua file with the Lua version of the language database.

```

11 local dbname = "language.dat.lua"
12 local language_dat
13 local dbfile = kpse.find_file(dbname)
14 if not dbfile then
15     err("file not found: "..dbname)
16 else
17     language_dat = dofile(dbfile)
18 end

```

Look up a language in the database, and return the associated information, as well as the canonical name of the language.

```

19 function lookupname(name)
20     if language_dat[name] then
21         return language_dat[name], name
22     else
23         for canon, data in pairs(language_dat) do
24             for _,syn in ipairs(data.synonyms) do
25                 if syn == name then
26                     return data, canon
27                 end
28             end
29         end
30     end
31 end

```

Set hyphenation patterns and exceptions for a language given by its name (in the database) and number (value of `\language`). Doesn't return anything, but will call `error()` if things go wrong.

```

32 function loadlanguage(lname, id)
33     local msg = "loading%s patterns and exceptions for: %s (\\language%d)"

```

Lookup the language in the database.

```

34     local ldata, cname = lookupname(lname)
35     if not ldata then
36         err("no entry in %s for this language: %s", dbname, lname)
37     end

```

Handle special languages.

```

38     if ldata.special then
39         if ldata.special == 'null' then
40             wlog(msg, ' (null)', cname, id)
41             return
42         elseif ldata.special:find('^disabled:') then

```

```

43         err("language disabled by %s: %s (%s)", dbname, cname,
44             ldata.special:gsub('^disabled:', ''))
45     elseif ldata.special == 0 then
46         err("\\language0 should be dumped in the format")
47     else
48         err("bad entry in %s for language %s")
49     end
50 end

```

The generic case: load hyphenation patterns and exceptions from files given by the language code.

```

51     wlog(msg, '', cname, id)
52     for ext, fun in pairs({pat = lang.patterns, hyp = lang.hyphenation}) do
53         local file = 'hyph-'..ldata.code..'..'..ext..'..txt'
54         local file = kpse.find_file(file) or err("file not found: %s", file)
55         local fh = io.open(file, 'r')
56         local data = fh:read('*a') or err("file not readable: %s", f)
57         fh:close()
58         fun(lang.new(id), data)
59     end
60 end
61 </lua>

```

2.2 hyphen.cfg

62 <*hyphen>

Start with unmodified code from babel.

```

63 \ifx\ProvidesFile\@undefined
64 \def\ProvidesFile#1[#2 #3 #4]{%
65 \wlog{File: #1 #4 #3 <#2>}%

```

Use a modified banner for LuaTeX.

```

66 \ifx\directlua\@undefined
67 \toks8{Babel <#3> and hyphenation patterns for }%
68 \else
69 \toks8{LuaTeX adaptation of babel <#3>
70 and hyphenation patterns for }%
71 \fi
72 \let\ProvidesFile\@undefined
73 }
74 \def\ProvidesLanguage#1[#2 #3 #4]{%
75 \wlog{Language: #1 #4 #3 <#2>}%
76 }
77 \else
78 \let\bbl@tempa\ProvidesFile
79 \def\ProvidesFile#1[#2 #3 #4]{%

```

Same here.

```

80 \ifx\directlua\@undefined
81 \toks8{Babel <#3> and hyphenation patterns for }%

```

```

82     \else
83         \toks8{LuaTeX adaptation of babel <#3>
84             and hyphenation patterns for }%
85     \fi
86     \bbl@tempa#1[#2 #3 #4]%
87     \let\ProvidesFile\bbl@tempa}
88 \def\ProvidesLanguage#1{%
89     \begingroup
90     \catcode'\ 10 %
91     \@makeother\/%
92     \@ifnextchar[%]
93         {\@provideslanguage{#1}}{\@provideslanguage{#1}[]}
94 \def\@provideslanguage#1[#2]{%
95     \wlog{Language: #1 #2}%
96     \expandafter\xdef\csname ver@#1.ldf\endcsname{#2}%
97     \endgroup}
98 \fi
99
    File identification is modified again.
100 \ProvidesFile{hyphen.cfg}
101         [2010/04/26 v3.81-luatex-1.3beta %
102         Language switching mechanism for LuaTeX, adapted from babel v3.81]
103 \ifx\AtBeginDocument\@undefined
104     \input plain.def\relax
105 \fi
106 \ifx\language\@undefined
107     \csname newcount\endcsname\language
108 \fi
109 \ifx\newlanguage\@undefined
110     \csname newcount\endcsname\last@language
111 \else
112     \countdef\last@language=19
113 \fi
114 \ifx\newlanguage\@undefined
115     \def\addlanguage#1{%
116         \global\advance\last@language \@ne
117         \ifnum\last@language<\@cclvi
118             \else
119                 \errmessage{No room for a new \string\language!}%
120             \fi
121         \global\chardef#1\last@language
122         \wlog{\string#1 = \string\language\the\last@language}}
123 \else
124     \def\addlanguage{\alloc@9\language\chardef\@cclvi}
125 \fi
126 \def\adddialect#1#2{%
127     \global\chardef#1#2\relax
128     \wlog{\string#1 = a dialect from \string\language#2}}
129 \def\iflanguage#1{%

```

```

130 \expandafter\ifx\csname l@#1\endcsname\relax
131 \@nolanerr{#1}%
132 \else
133 \bbl@afterfi{\ifnum\csname l@#1\endcsname=\language
134 \expandafter\@firstoftwo
135 \else
136 \expandafter\@secondoftwo
137 \fi}%
138 \fi}
139 \edef\selectlanguage{%
140 \noexpand\protect
141 \expandafter\noexpand\csname selectlanguage \endcsname
142 }
143 \ifx\@undefined\protect\let\protect\relax\fi
144 \ifx\documentclass\@undefined
145 \def\xstring{\string\string\string}
146 \else
147 \let\xstring\string
148 \fi
149 \xdef\bbl@language@stack{}
150 \def\bbl@push@language{%
151 \xdef\bbl@language@stack{\language+\bbl@language@stack}%
152 }
153 \def\bbl@pop@lang#1+#2-#3{%
154 \def\language{#1}\xdef#3{#2}%
155 }
156 \def\bbl@pop@language{%
157 \expandafter\bbl@pop@lang\bbl@language@stack-\bbl@language@stack
158 \expandafter\bbl@set@language\expandafter{\language}%
159 }
160 \expandafter\def\csname selectlanguage \endcsname#1{%
161 \bbl@push@language
162 \aftergroup\bbl@pop@language
163 \bbl@set@language{#1}}
164 \def\bbl@set@language#1{%
165 \edef\language{#1}
166 \ifnum\escapechar=\expandafter'\string#1\@empty
167 \else \string#1\@empty\fi}%
168 \select@language{\language}%
169 \if@filesw
170 \protected@write\@auxout{}\{\string\select@language{\language}\}%
171 \addtocontents{toc}\{\xstring\select@language{\language}\}%
172 \addtocontents{lof}\{\xstring\select@language{\language}\}%
173 \addtocontents{lot}\{\xstring\select@language{\language}\}%
174 \fi}
175 \def\select@language#1{%
176 \expandafter\ifx\csname l@#1\endcsname\relax
177 \@nolanerr{#1}%
178 \else
179 \expandafter\ifx\csname date#1\endcsname\relax

```

```

180     \@noopterr{#1}%
181   \else
182     \bbl@patterns{\language}%
183     \originalTeX
184     \expandafter\def\expandafter\originalTeX
185       \expandafter{\csname noextras#1\endcsname
186         \let\originalTeX\@empty}%
187     \languageshorthands{none}%
188     \babel@beginsave
189     \csname captions#1\endcsname
190     \csname date#1\endcsname
191     \csname extras#1\endcsname\relax
192     \babel@savevariable\lefthyphenmin
193     \babel@savevariable\righthyphenmin
194     \expandafter\ifx\csname #1hyphenmins\endcsname\relax
195       \set@hyphenmins\tw@\thr@@\relax
196   \else
197     \expandafter\expandafter\expandafter\set@hyphenmins
198       \csname #1hyphenmins\endcsname\relax
199   \fi
200 \fi
201 \fi}
202 \long\def\otherlanguage#1{%
203   \csname selectlanguage \endcsname{#1}%
204   \ignorespaces
205 }
206 \long\def\endotherlanguage{%
207   \originalTeX
208   \global\@ignoretrue\ignorespaces
209 }
210 \expandafter\def\csname otherlanguage*\endcsname#1{%
211   \foreign@language{#1}%
212 }
213 \expandafter\def\csname endotherlanguage*\endcsname{%
214   \csname noextras\language\endcsname
215 }
216 \def\foreignlanguage{\protect\csname foreignlanguage \endcsname}
217 \expandafter\def\csname foreignlanguage \endcsname#1#2{%
218   \begingroup
219     \originalTeX
220     \foreign@language{#1}%
221     #2%
222     \csname noextras#1\endcsname
223   \endgroup
224 }
225 \def\foreign@language#1{%
226   \def\language{#1}%
227   \expandafter\ifx\csname l@#1\endcsname\relax
228     \@nolanerr{#1}%
229   \else

```

```

230 \bbl@patterns{\language}%
231 \languageshorthands{none}%
232 \csname extras#1\endcsname
233 \expandafter\ifx\csname #1hyphenmins\endcsname\relax
234 \set@hyphenmins\tw@\thr@@\relax
235 \else
236 \expandafter\expandafter\expandafter\set@hyphenmins
237 \csname #1hyphenmins\endcsname\relax
238 \fi
239 \fi
240 }
241 \def\bbl@patterns#1{%
242 \language=\expandafter\ifx\csname l@#1:\f@encoding\endcsname\relax
243 \csname l@#1\endcsname
244 \else
245 \csname l@#1:\f@encoding\endcsname
246 \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.

```

247 \ifx\directlua@\undefined\else
248 \ifx\directlua\relax\else
249 \ifcsname lu@texhyphen@loaded@\the\language\endcsname \else
250 \global\@namedef{lu@texhyphen@loaded@\the\language}{}%
251 \directlua{
252     if not luatexhyphen then
253         dofile(assert(kpse.find_file("luatex-hyphen.lua")))
254     end
255     luatexhyphen.loadlanguage("\luatexluaescapestring{#1}",
256 \the\language)}%
257 \fi
258 \fi
259 \fi
260 }
261 \def\hyphenrules#1{%
262 \expandafter\ifx\csname l@#1\endcsname@\undefined
263 \@nolanerr{#1}%
264 \else
265 \bbl@patterns{#1}%
266 \languageshorthands{none}%
267 \expandafter\ifx\csname #1hyphenmins\endcsname\relax
268 \set@hyphenmins\tw@\thr@@\relax
269 \else
270 \expandafter\expandafter\expandafter\set@hyphenmins
271 \csname #1hyphenmins\endcsname\relax
272 \fi
273 \fi

```



```

274 }
275 \def\endhyphenrules{}
276 \def\providehyphenmins#1#2{%
277   \expandafter\ifx\csname #1hyphenmins\endcsname\relax
278     \@namedef{#1hyphenmins}{#2}%
279   \fi}
280 \def\set@hyphenmins#1#2{\lefthyphenmin#1\righthyphenmin#2}
281 \def\LdfInit{%
282   \chardef\atcatcode=\catcode'\@
283   \catcode'\@=11\relax
284   \input babel.def\relax
285   \catcode'\@=\atcatcode \let\atcatcode\relax
286   \LdfInit}
287 \ifx\originalTeX\@undefined\let\originalTeX\@empty\fi
288 \ifx\babel@beginsave\@undefined\let\babel@beginsave\relax\fi
289 \ifx\PackageError\@undefined
290   \def\@nolanerr#1{%
291     \errhelp{Your command will be ignored, type <return> to proceed}%
292     \errmessage{You haven't defined the language #1\space yet}}
293   \def\@nopatterns#1{%
294     \message{No hyphenation patterns were loaded for}%
295     \message{the language '#1'}%
296     \message{I will use the patterns loaded for \string\language=0
297       instead}}
298   \def\@noopterr#1{%
299     \errmessage{The option #1 was not specified in \string\usepackage}
300     \errhelp{You may continue, but expect unexpected results}}
301   \def\@activated#1{%
302     \wlog{Package babel Info: Making #1 an active character}}
303 \else
304   \newcommand*\@nolanerr[1]{%
305     \PackageError{babel}%
306       {You haven't defined the language #1\space yet}%
307       {Your command will be ignored, type <return> to proceed}}
308   \newcommand*\@nopatterns[1]{%
309     \PackageWarningNoLine{babel}%
310       {No hyphenation patterns were loaded for\MessageBreak
311         the language '#1'\MessageBreak
312         I will use the patterns loaded for \string\language=0
313         instead}}
314   \newcommand*\@noopterr[1]{%
315     \PackageError{babel}%
316       {You haven't loaded the option #1\space yet}%
317       {You may proceed, but expect unexpected results}}
318   \newcommand*\@activated[1]{%
319     \PackageInfo{babel}{%
320       Making #1 an active character}}
321 \fi
322 \def\process@line#1#2 #3/{%
323   \ifx=#1

```

```

324 \process@synonym#2 /
325 \else
326 \process@language#1#2 #3/%
327 \fi
328 }
329 \toks@{}
330 \def\process@synonym#1 /{%
331 \ifnum\last@language=\m@ne
332 \expandafter\chardef\csname l@#1\endcsname\relax
333 \wlog{\string\l@#1=\string\language0}
334 \toks@\expandafter{\the\toks@
335 \expandafter\let\csname #1hyphenmins\expandafter\endcsname
336 \csname\language\hyphenmins\endcsname}%
337 \else
338 \expandafter\chardef\csname l@#1\endcsname\last@language
339 \wlog{\string\l@#1=\string\language\the\last@language}
340 \expandafter\let\csname #1hyphenmins\expandafter\endcsname
341 \csname\language\hyphenmins\endcsname
342 \fi
343 }
344 \def\process@language#1 #2 #3/{%
345 \expandafter\addlanguage\csname l@#1\endcsname
346 \expandafter\language\csname l@#1\endcsname
347 \def\language{#1}%

```

In the Lua_{TEX} case, we have to decide whether to load the language now. Remember our choice, since we'll need it two times more.

If we choose to load the language now, mark it as loaded. This is done using _{TEX} macros in order to survive the format dumping-loading cycle, which would not be as straightforward using Lua objects.

```

348 \ifx\directlua@\undefined
349 \global\toks8\expandafter{\the\toks8#1, }%
350 \else
351 \directlua{
352   if not luatexhyphen then
353     dofile(assert(kpse.find_file("luatex-hyphen.lua")))
354   end
355   processnow = (tex.language == 0) or
356     (luatexhyphen.lookupname("\luatexluaescapestring{#1}") == nil)}%
357 \ifnum0=\directlua{tex.sprint(processnow and "0" or "1")}\relax
358 \global\toks8\expandafter{\the\toks8#1, }%
359 \global\@namedef{lu@texhyphen@loaded@the\language}{}%
360 \fi
361 \fi
362 \begingroup
363 \bbl@get@enc#1:\@@@
364 \ifx\bbl@hyph@enc\@empty
365 \else
366 \fontencoding{\bbl@hyph@enc}\selectfont
367 \fi

```

```

368 \lefthyphenmin\m@ne
    Conditionally input the patterns file.
369 \ifx\directlua\@undefined
370 \input #2\relax
371 \else
372 \ifnum0=\directlua{tex.sprint(processnow and "0" or "1")}\relax
373 \input #2\relax
374 \fi
375 \fi
376 \ifnum\lefthyphenmin=\m@ne
377 \else
378 \expandafter\xdef\csname #1hyphenmins\endcsname{%
379 \the\lefthyphenmin\the\righthyphenmin}%
380 \fi
381 \endgroup
382 \ifnum\the\language=\z@
383 \expandafter\ifx\csname #1hyphenmins\endcsname\relax
384 \set@hyphenmins\tw@\thr@@\relax
385 \else
386 \expandafter\expandafter\expandafter\set@hyphenmins
387 \csname #1hyphenmins\endcsname
388 \fi
389 \the\toks@
390 \fi
391 \toks@{}%
392 \def\bbl@tempa{#3}%
393 \ifx\bbl@tempa\@empty
394 \else
395 \ifx\bbl@tempa\space
396 \else
    Conditionnaly input the exceptions file.
397 \ifx\directlua\@undefined
398 \input #3\relax
399 \else
400 \ifnum0=\directlua{tex.sprint(processnow and "0" or "1")}\relax
401 \input #3\relax
402 \fi
403 \directlua{processnow = nil}%
404 \fi
405 \fi
406 \fi
407 }
408 \def\bbl@get@enc#1:#2\@@@{%
409 \def\bbl@tempa{#1}%
410 \def\bbl@tempb{#2}%
411 \ifx\bbl@tempb\@empty
412 \let\bbl@hyph@enc\@empty
413 \else

```

```

414 \bbl@get@enc#2\@@@
415 \edef\bbl@hyph@enc{\bbl@tempa}%
416 \fi}
417 \openin1 = language.dat
418 \ifeof1
419 \message{I couldn't find the file language.dat,\space
420         I will try the file hyphen.tex}
421 \input hyphen.tex\relax
422 \else
423 \last@language@m@ne
424 \loop
425 \endlinechar@m@ne
426 \read1 to \bbl@line
427 \endlinechar'\^^M
428 \ifx\bbl@line\@empty
429 \else
430 \edef\bbl@line{\bbl@line\space/}%
431 \expandafter\process@line\bbl@line
432 \fi
433 \iftrue \csname fi\endcsname
434 \csname if\ifeof1 false\else true\fi\endcsname
435 \repeat
436 \language=0
437 \fi
438 \closein1
439 \let\process@language\@undefined
440 \let\process@synonym\@undefined
441 \let\process@line\@undefined
442 \let\bbl@tempa\@undefined
443 \let\bbl@tempb\@undefined
444 \let\bbl@eq\@undefined
445 \let\bbl@line\@undefined
446 \let\bbl@get@enc\@undefined
447 \ifx\addto@hook\@undefined
448 \else
449 \expandafter\addto@hook\expandafter\everyjob\expandafter{%
450 \expandafter\typeout\expandafter{\the\toks8 loaded.}}
451 \fi
452 \</hyphen>

```