# The `hyphen.cfg` file for LuaTeX

Khaled Hosny, Élie Roux, and Manuel Pégourié-Gonnard

`khaledhosny@eglug.org`

`elie.roux@telecom-bretagne.eu`

`mpg@elzevir.fr`

2012/04/16 v1.5

**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 LuaTeX'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.

There is a version of `etex.src` modified for the same reasons using similar code, which also makes use of the `luatex-hyphen.lua` and `language.dat.lua` files described here.

## 1 Documentation

Hyphenation patterns should be loaded at runtime with LuaTeX: 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 LuaTeX, 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 one mandatory field:

   ```
   synonyms = { <string> alternative name, ...}
   ```

This field's value must be the same as in `language.dat`.

2. A variable part consisting of either:

   - For most languages:

     ```
     patterns = <string> filenames for patterns
     hyphenation = <string> filenames for exceptions
     ```

     Each string contains a coma-separated list of file names (whitespace before or after the coma is not accepted). The files given by `patterns` (resp. `hypenation`) must be plain text files encoded in UTF-8, with only patterns (resp. exceptions) and not even comments: their content will be used directly without being parsed by TeX. If one of these keys is missing or is the empty string, it is ignored and no patterns (resp. exceptions) are loaded for this language.

   - Special cases are supported by a field `special`. Currently, the following kind of values are recognized:

     `'disabled:<reason>'` allows to disable specific languages: when the user tries to load this language, an error will be issued, with the `<reason>`.

     `'language0'` only `english` should use this type of special, to indicate it is normally dumped in the format as `\language0` (see below).

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

3. Optional fields may be added. For example:

   ```
   loader = <string> name of the TeX loader
   lefthyphenmin = <number> value for \lefthyphenmin
   righthyphenmin = <number> value for \righthyphenmin
   ```

   Those fields are present in `language.dat.lua` as generated by `tlmgr`, for example, but they *are not* used by the present code in any way.

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.

---

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

# 2 Implementation

## 2.1 luatex-hyphen.lua

1 ⟨∗lua⟩

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

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

Two functions for error and information reporting.

```lua
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.

```lua
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     wlog('using data file: %s', dbfile)
18     language_dat = dofile(dbfile)
19 end
```

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

```lua
20 function lookupname(name)
21     if language_dat[name] then
22         return language_dat[name], name
23     else
24         for canon, data in pairs(language_dat) do
25             for _,syn in ipairs(data.synonyms) do
26                 if syn == name then
27                     return data, canon
28                 end
29             end
30         end
31     end
32 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.

```lua
33 function loadlanguage(lname, id)
34     if id == 0 then
35         return
```

```
36      end
37      local msg = "loading%s patterns and exceptions for: %s (\\language%d)"
```

Lookup the language in the database.

```
38      local ldata, cname = lookupname(lname)
39      if not ldata then
40          err("no entry in %s for this language: %s", dbname, lname)
41      end
```

Handle special languages.

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

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

```
52      wlog(msg, '', cname, id)
53      for _, item in ipairs{'patterns', 'hyphenation'} do
54          local filelist = ldata[item]
55          if filelist ~= nil and filelist ~= '' then
56            for _, file in ipairs(filelist:explode(',')) do
57              local file = kpse.find_file(file) or err("file not found: %s", file)
58              local fh = io.open(file, 'r')
59              local data = fh:read('*a') or err("file not readable: %s", f)
60              fh:close()
61              lang[item](lang.new(id), data)
62            end
63          else
64              if item == 'hyphenation' then item = item..' exceptions' end
65              wlog("info: no %s for this language", item)
66          end
67      end
68 end
69 function adddialect(dialect, language)
70      if dialect ~= '0' then
71          dialect = dialect:gsub('l@', '')
72          language = language:gsub('l@', '')
73          data = language_dat[language]
74          if data then
75              data.synonyms[#data.synonyms+1] = dialect
76          end
77      end
78 end

79 ⟨/lua⟩
```

4

## 2.2   hyphen.cfg

80 ⟨∗hyphen⟩

   Start with unmodified code from babel.

```
81 \ifx\ProvidesFile\@undefined
82   \def\ProvidesFile#1[#2 #3 #4]{%
83     \wlog{File: #1 #4 #3 <#2>}%
```

Use a modified banner for LuaTEX.

```
84     \ifx\directlua\@undefined
85       \toks8{Babel <#3> and hyphenation patterns for }%
86     \else
87       \toks8{LuaTeX adaptation of babel <#3>
88         and hyphenation patterns for }%
89     \fi

90     \let\ProvidesFile\@undefined
91     }
92   \def\ProvidesLanguage#1[#2 #3 #4]{%
93     \wlog{Language: #1 #4 #3 <#2>}%
94     }
95 \else
96   \let\bbl@tempa\ProvidesFile
97   \def\ProvidesFile#1[#2 #3 #4]{%
```

Same here.

```
98     \ifx\directlua\@undefined
99       \toks8{Babel <#3> and hyphenation patterns for }%
100    \else
101      \toks8{LuaTeX adaptation of babel <#3>
102        and hyphenation patterns for }%
103    \fi

104    \bbl@tempa#1[#2 #3 #4]%
105    \let\ProvidesFile\bbl@tempa}
106  \def\ProvidesLanguage#1{%
107    \begingroup
108      \catcode`\ 10 %
109      \@makeother\/%
110      \@ifnextchar[%
111        {\@provideslanguage{#1}}{\@provideslanguage{#1}[]}}
112  \def\@provideslanguage#1[#2]{%
113    \wlog{Language: #1 #2}%
114    \expandafter\xdef\csname ver@#1.ldf\endcsname{#2}%
115    \endgroup}
116 \fi
117
```

   File identification is modified again.

```
118 \ProvidesFile{hyphen.cfg}
119              [2012/04/16 v3.8l-luatex-1.5 %
120       Language switching mechanism for LuaTeX, adapted from babel v3.8m]
```

```
121 \ifx\AtBeginDocument\@undefined
122   \input plain.def\relax
123 \fi
124 \ifx\language\@undefined
125   \csname newcount\endcsname\language
126 \fi
127 \ifx\newlanguage\@undefined
128   \csname newcount\endcsname\last@language
129 \else
130   \countdef\last@language=19
131 \fi
132 \ifx\newlanguage\@undefined
133   \def\addlanguage#1{%
134     \global\advance\last@language \@ne
135     \ifnum\last@language<\@cclvi
136     \else
137       \errmessage{No room for a new \string\language!}%
138     \fi
139     \global\chardef#1\last@language
140     \wlog{\string#1 = \string\language\the\last@language}}
141 \else
142   \def\addlanguage{\alloc@9\language\chardef\@cclvi}
143 \fi
144 \def\adddialect#1#2{%
145     \global\chardef#1#2\relax
146     \ifx\directlua\@undefined\else
147       \ifx\directlua\relax\else
148         \directlua{
149           if not luatexhyphen then
150             dofile(assert(kpse.find_file("luatex-hyphen.lua")))
151           end
152           luatexhyphen.adddialect("\string#1", "\string#2")
153         }%
154       \fi
155     \fi
156     \wlog{\string#1 = a dialect from \string\language#2}}
157 \def\iflanguage#1{%
158   \expandafter\ifx\csname l@#1\endcsname\relax
159     \@nolanerr{#1}%
160   \else
161     \bbl@afterfi{\ifnum\csname l@#1\endcsname=\language
162       \expandafter\@firstoftwo
163     \else
164       \expandafter\@secondoftwo
165     \fi}%
166   \fi}
167 \edef\selectlanguage{%
168   \noexpand\protect
169   \expandafter\noexpand\csname selectlanguage \endcsname
170   }
```

```
171 \ifx\@undefined\protect\let\protect\relax\fi
172 \ifx\documentclass\@undefined
173   \def\xstring{\string\string\string}
174 \else
175   \let\xstring\string
176 \fi
177 \xdef\bbl@language@stack{}
178 \def\bbl@push@language{%
179   \xdef\bbl@language@stack{\languagename+\bbl@language@stack}%
180   }
181 \def\bbl@pop@lang#1+#2-#3{%
182   \def\languagename{#1}\xdef#3{#2}%
183   }
184 \def\bbl@pop@language{%
185   \expandafter\bbl@pop@lang\bbl@language@stack-\bbl@language@stack
186   \expandafter\bbl@set@language\expandafter{\languagename}%
187   }
188 \expandafter\def\csname selectlanguage \endcsname#1{%
189   \bbl@push@language
190   \aftergroup\bbl@pop@language
191   \bbl@set@language{#1}}
192 \def\bbl@set@language#1{%
193   \edef\languagename{%
194     \ifnum\escapechar=\expandafter`\string#1\@empty
195     \else \string#1\@empty\fi}%
196   \select@language{\languagename}%
197   \if@filesw
198     \protected@write\@auxout{}{\string\select@language{\languagename}}%
199     \addtocontents{toc}{\xstring\select@language{\languagename}}%
200     \addtocontents{lof}{\xstring\select@language{\languagename}}%
201     \addtocontents{lot}{\xstring\select@language{\languagename}}%
202   \fi}
203 \def\select@language#1{%
204   \expandafter\ifx\csname l@#1\endcsname\relax
205     \@nolanerr{#1}%
206   \else
207     \expandafter\ifx\csname date#1\endcsname\relax
208       \@noopterr{#1}%
209     \else
210       \bbl@patterns{\languagename}%
211       \originalTeX
212       \expandafter\def\expandafter\originalTeX
213         \expandafter{\csname noextras#1\endcsname
214                     \let\originalTeX\@empty}%
215       \languageshorthands{none}%
216       \babel@beginsave
217       \csname captions#1\endcsname
218       \csname date#1\endcsname
219       \csname extras#1\endcsname\relax
220       \babel@savevariable\lefthyphenmin
```

```
221      \babel@savevariable\righthyphenmin
222      \expandafter\ifx\csname #1hyphenmins\endcsname\relax
223        \set@hyphenmins\tw@\thr@@\relax
224      \else
225        \expandafter\expandafter\expandafter\set@hyphenmins
226          \csname #1hyphenmins\endcsname\relax
227      \fi
228    \fi
229  \fi}
230 \long\def\otherlanguage#1{%
231   \csname selectlanguage \endcsname{#1}%
232   \ignorespaces
233   }
234 \long\def\endotherlanguage{%
235   \originalTeX
236   \global\@ignoretrue\ignorespaces
237   }
238 \expandafter\def\csname otherlanguage*\endcsname#1{%
239   \foreign@language{#1}%
240   }
241 \expandafter\def\csname endotherlanguage*\endcsname{%
242   \csname noextras\languagename\endcsname
243   }
244 \def\foreignlanguage{\protect\csname foreignlanguage \endcsname}
245 \expandafter\def\csname foreignlanguage \endcsname#1#2{%
246   \begingroup
247     \originalTeX
248     \foreign@language{#1}%
249     #2%
250     \csname noextras#1\endcsname
251   \endgroup
252   }
253 \def\foreign@language#1{%
254   \def\languagename{#1}%
255   \expandafter\ifx\csname l@#1\endcsname\relax
256     \@nolanerr{#1}%
257   \else
258     \bbl@patterns{\languagename}%
259     \languageshorthands{none}%
260     \csname extras#1\endcsname
261     \expandafter\ifx\csname #1hyphenmins\endcsname\relax
262       \set@hyphenmins\tw@\thr@@\relax
263     \else
264       \expandafter\expandafter\expandafter\set@hyphenmins
265         \csname #1hyphenmins\endcsname\relax
266     \fi
267   \fi
268   }
269 \def\bbl@patterns#1{%
270   \language=\expandafter\ifx\csname l@#1:\f@encoding\endcsname\relax
```

```
271    \csname l@#1\endcsname
272   \else
273    \csname l@#1:\f@encoding\endcsname
274   \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.

```
275   \ifx\directlua\@undefined\else
276     \ifx\directlua\relax\else
277       \ifcsname lu@texhyphen@loaded@\the\language\endcsname \else
278         \global\@namedef{lu@texhyphen@loaded@\the\language}{}%
279         \directlua{
280           if not luatexhyphen then
281              dofile(assert(kpse.find_file("luatex-hyphen.lua")))
282           end
283           luatexhyphen.loadlanguage("\luatexluaescapestring{#1}",
284             \the\language)}%
285       \fi
286     \fi
287   \fi
288 }
289 \def\hyphenrules#1{%
290   \expandafter\ifx\csname l@#1\endcsname\@undefined
291     \@nolanerr{#1}%
292   \else
293     \bbl@patterns{#1}%
294     \languageshorthands{none}%
295       \expandafter\ifx\csname #1hyphenmins\endcsname\relax
296         \set@hyphenmins\tw@\thr@@\relax
297       \else
298         \expandafter\expandafter\expandafter\set@hyphenmins
299         \csname #1hyphenmins\endcsname\relax
300       \fi
301   \fi
302   }
303 \def\endhyphenrules{}
304 \def\providehyphenmins#1#2{%
305   \expandafter\ifx\csname #1hyphenmins\endcsname\relax
306     \@namedef{#1hyphenmins}{#2}%
307   \fi}
308 \def\set@hyphenmins#1#2{\lefthyphenmin#1\righthyphenmin#2}
309 \def\LdfInit{%
310   \chardef\atcatcode=\catcode`\@
311   \catcode`\@=11\relax
312   \input babel.def\relax
313   \catcode`\@=\atcatcode \let\atcatcode\relax
314   \LdfInit}
```

```
315 \ifx\originalTeX\@undefined\let\originalTeX\@empty\fi
316 \ifx\babel@beginsave\@undefined\let\babel@beginsave\relax\fi
317 \ifx\PackageError\@undefined
318   \def\@nolanerr#1{%
319     \errhelp{Your command will be ignored, type <return> to proceed}%
320     \errmessage{You haven't defined the language #1\space yet}}
321   \def\@nopatterns#1{%
322     \message{No hyphenation patterns were loaded for}%
323     \message{the language '#1'}%
324     \message{I will use the patterns loaded for \string\language=0
325          instead}}
326   \def\@noopterr#1{%
327     \errmessage{The option #1 was not specified in \string\usepackage}
328     \errhelp{You may continue, but expect unexpected results}}
329   \def\@activated#1{%
330     \wlog{Package babel Info: Making #1 an active character}}
331 \else
332   \newcommand*{\@nolanerr}[1]{%
333     \PackageError{babel}%
334                  {You haven't defined the language #1\space yet}%
335         {Your command will be ignored, type <return> to proceed}}
336   \newcommand*{\@nopatterns}[1]{%
337     \PackageWarningNoLine{babel}%
338         {No hyphenation patterns were loaded for\MessageBreak
339           the language '#1'\MessageBreak
340           I will use the patterns loaded for \string\language=0
341           instead}}
342   \newcommand*{\@noopterr}[1]{%
343     \PackageError{babel}%
344                  {You haven't loaded the option #1\space yet}%
345             {You may proceed, but expect unexpected results}}
346   \newcommand*{\@activated}[1]{%
347     \PackageInfo{babel}{%
348       Making #1 an active character}}
349 \fi
350 \def\process@line#1#2 #3/{%
351   \ifx=#1
352     \process@synonym#2 /
353   \else
354     \process@language#1#2 #3/%
355   \fi
356   }
357 \toks@{}
358 \def\process@synonym#1 /{%
359   \ifnum\last@language=\m@ne
360     \expandafter\chardef\csname l@#1\endcsname0\relax
361     \wlog{\string\l@#1=\string\language0}
362     \toks@\expandafter{\the\toks@
363       \expandafter\let\csname #1hyphenmins\expandafter\endcsname
364       \csname\languagename hyphenmins\endcsname}%
```

10

```
365    \else
366      \expandafter\chardef\csname l@#1\endcsname\last@language
367      \wlog{\string\l@#1=\string\language\the\last@language}
368      \expandafter\let\csname #1hyphenmins\expandafter\endcsname
369      \csname\languagename hyphenmins\endcsname
370    \fi
371    }
372 \def\process@language#1 #2 #3/{%
373    \expandafter\addlanguage\csname l@#1\endcsname
374    \expandafter\language\csname l@#1\endcsname
375    \def\languagename{#1}%
```

In the LuaTeX case, we have to decide wether 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 straigthforward using Lua objects.

```
376    \ifx\directlua\@undefined
377      \global\toks8\expandafter{\the\toks8#1, }%
378    \else
379      \directlua{
380        if not luatexhyphen then
381          dofile(assert(kpse.find_file("luatex-hyphen.lua")))
382        end
383        processnow = (tex.language == 0) or
384          (luatexhyphen.lookupname("\luatexluaescapestring{#1}") == nil)}%
385      \ifnum0=\directlua{tex.sprint(processnow and "0" or "1")}\relax
386        \global\toks8\expandafter{\the\toks8#1, }%
387        \global\@namedef{lu@texhyphen@loaded@\the\language}{}%
388      \fi
389    \fi

390    \begingroup
391      \bbl@get@enc#1:\@@@
392      \ifx\bbl@hyph@enc\@empty
393      \else
394        \fontencoding{\bbl@hyph@enc}\selectfont
395      \fi
396      \lefthyphenmin\m@ne
```

Conditionally input the patterns file.

```
397      \ifx\directlua\@undefined
398        \input #2\relax
399      \else
400        \ifnum0=\directlua{tex.sprint(processnow and "0" or "1")}\relax
401          \input #2\relax
402        \fi
403      \fi

404      \ifnum\lefthyphenmin=\m@ne
405      \else
406        \expandafter\xdef\csname #1hyphenmins\endcsname{%
```

```
407        \the\lefthyphenmin\the\righthyphenmin}%
408      \fi
409    \endgroup
410    \ifnum\the\language=\z@
411      \expandafter\ifx\csname #1hyphenmins\endcsname\relax
412        \set@hyphenmins\tw@\thr@@\relax
413      \else
414        \expandafter\expandafter\expandafter\set@hyphenmins
415          \csname #1hyphenmins\endcsname
416      \fi
417      \the\toks@
418    \fi
419    \toks@{}%
420    \def\bbl@tempa{#3}%
421    \ifx\bbl@tempa\@empty
422    \else
423      \ifx\bbl@tempa\space
424      \else
```

Conditionnaly input the exceptions file.

```
425        \ifx\directlua\@undefined
426          \input #3\relax
427        \else
428          \ifnum0=\directlua{tex.sprint(processnow and "0" or "1")}\relax
429            \input #3\relax
430          \fi
431          \directlua{processnow = nil}%
432        \fi

433      \fi
434    \fi
435    }
436 \def\bbl@get@enc#1:#2\@@@{%
437    \def\bbl@tempa{#1}%
438    \def\bbl@tempb{#2}%
439    \ifx\bbl@tempb\@empty
440      \let\bbl@hyph@enc\@empty
441    \else
442      \bbl@get@enc#2\@@@
443      \edef\bbl@hyph@enc{\bbl@tempa}%
444    \fi}
445 \openin1 = language.dat
446 \ifeof1
447    \message{I couldn't find the file language.dat,\space
448            I will try the file hyphen.tex}
449    \input hyphen.tex\relax
450    \def\l@english{0}%
451    \def\languagename{english}%
452 \else
453    \last@language\m@ne
454    \loop
```

```
455    \endlinechar\m@ne
456    \read1 to \bbl@line
457    \endlinechar`\^^M
458    \ifx\bbl@line\@empty
459    \else
460      \edef\bbl@line{\bbl@line\space/}%
461      \expandafter\process@line\bbl@line
462      \ifx\bbl@defaultlanguage\@undefined
463        \let\bbl@defaultlanguage\languagename
464      \fi
465    \fi
466    \iftrue \csname fi\endcsname
467    \csname if\ifeof1 false\else true\fi\endcsname
468  \repeat
469  \language=0
470  \let\languagename\bbl@defaultlanguage
471  \let\bbl@defaultlanguage\@undefined
472 \fi
473 \closein1
474 \let\process@language\@undefined
475 \let\process@synonym\@undefined
476 \let\process@line\@undefined
477 \let\bbl@tempa\@undefined
478 \let\bbl@tempb\@undefined
479 \let\bbl@eq@\@undefined
480 \let\bbl@line\@undefined
481 \let\bbl@get@enc\@undefined
482 \ifx\addto@hook\@undefined
483 \else
484   \expandafter\addto@hook\expandafter\everyjob\expandafter{%
485     \expandafter\typeout\expandafter{\the\toks8 loaded.}}
486 \fi

487 ⟨/hyphen⟩
```