

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`.

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 wlog(msg, ...)
9     texio.write_nl('log', 'luatex-hyphen: '..string.format(msg, ...))
10 end
11
12 local function err(msg, ...)
13     error('luatex-hyphen: '..string.format(msg, ...))
14 end
15
16 luatexhyphen.language_dat = {}
17 local dbfile = kpse.find_file(dbname)
18 if not dbfile then
19     err("file not found: "..dbname)
20 else
21     luatexhyphen.language_dat = dofile(dbfile)
22 end
23
24 function luatexhyphen.lookupname(l)
25     if luatexhyphen.language_dat[l] then
26         return luatexhyphen.language_dat[l], 1
27     else
28         for orig,lt in pairs(luatexhyphen.language_dat) do
29             for _,syn in ipairs(lt.synonyms) do
30                 if syn == l then
31                     return lt, orig
32                 end
33             end
34         end
35     end
36     return nil
37 end
38
39 function luatexhyphen.loadlanguage(l, id)
40     local lt, orig = luatexhyphen.lookupname(l)
41     if not lt then
42         err("no entry in %s for this language: %s", dbname, l)
43         return
44     end
45     local msg = "loading%s patterns and exceptions for: %s (\\language%d)"
46     if lt.special then
47         if lt.special == 'null' then
48             wlog(msg, ' (null)', orig, id)
49         elseif lt.special:find('^disabled:') then
50             err("language disabled by %s: %s (%s)", dbname, orig,
51                 lt.special:gsub('^disabled:', ''))
52         elseif lt.special == 0 then
53             err("\\language0 should be dumped in the format")
54         else
55             err("bad entry in %s for language %s")
56         end
57     end
58     return

```

```

58     end
59     wlog(msg, '', orig, id)
60     for ext, fun in pairs({pat = lang.patterns, hyp = lang.hyphenation}) do
61         local n = 'hyph-'..lt.code..'..'..ext..'..txt'
62         local f = kpse.find_file(n)
63         if not f then
64             err("file not found: %s", n)
65             return
66         end
67         f = io.open(f, 'r')
68         local data = f:read('*a')
69         f:close()
70         if not data then
71             err("file not readable: %s", f)
72             return
73         end
74         fun(lang.new(id), data)
75     end
76 end

```

2.2 hyphen.cfg

```

77 \ifx\ProvidesFile\@undefined
78 \def\ProvidesFile#1[#2 #3 #4]{%
79 \wlog{File: #1 #4 #3 <#2>}%

    Use a modified banner for LuaTeX.
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 \let\ProvidesFile\@undefined
87 }
88 \def\ProvidesLanguage#1[#2 #3 #4]{%
89 \wlog{Language: #1 #4 #3 <#2>}%
90 }
91 \else
92 \let\bbl@tempa\ProvidesFile
93 \def\ProvidesFile#1[#2 #3 #4]{%

    Same here.
94 \ifx\directlua\@undefined
95 \toks8{Babel <#3> and hyphenation patterns for }%
96 \else
97 \toks8{LuaTeX adaptation of babel <#3>
98 and hyphenation patterns for }%
99 \fi
100 \bbl@tempa#1[#2 #3 #4]%
101 \let\ProvidesFile\bbl@tempa}

```

```

102 \def\ProvidesLanguage#1{%
103   \begingroup
104   \catcode'\ 10 %
105   \@makeother\/%
106   \@ifnextchar[%]
107     {\@provideslanguage{#1}}{\@provideslanguage{#1}[]}}
108 \def\@provideslanguage#1[#2]{%
109   \wlog{Language: #1 #2}%
110   \expandafter\xdef\csname ver@#1.ldf\endcsname{#2}%
111   \endgroup}
112 \fi
113
114 \ProvidesFile{hyphen.cfg}
115   [2010/04/26 v3.81-luatex-1.3beta %
116   Language switching mechanism for LuaTeX, adapted from babel v3.81]
117 \ifx\AtBeginDocument\@undefined
118   \input plain.def\relax
119 \fi
120 \ifx\language\@undefined
121   \csname newcount\endcsname\language
122 \fi
123 \ifx\newlanguage\@undefined
124   \csname newcount\endcsname\last@language
125 \else
126   \countdef\last@language=19
127 \fi
128 \ifx\newlanguage\@undefined
129   \def\addlanguage#1{%
130     \global\advance\last@language \@ne
131     \ifnum\last@language<\ccclvi
132       \else
133         \errmessage{No room for a new \string\language!}%
134       \fi
135     \global\chardef#1\last@language
136     \wlog{\string#1 = \string\language\the\last@language}}
137 \else
138   \def\addlanguage{\alloc@9\language\chardef\ccclvi}
139 \fi
140 \def\adddialect#1#2{%
141   \global\chardef#1#2\relax
142   \wlog{\string#1 = a dialect from \string\language#2}}
143 \def\iflanguage#1{%
144   \expandafter\ifx\csname l@#1\endcsname\relax
145     \@nolanerr{#1}%
146   \else
147     \bbl@afterfi{\ifnum\csname l@#1\endcsname=\language
148       \expandafter\@firstoftwo
149     \else

```

```

150     \expandafter\@secondoftwo
151     \fi}%
152 \fi}
153 \edef\selectlanguage{%
154   \noexpand\protect
155   \expandafter\noexpand\csname selectlanguage \endcsname
156 }
157 \ifx\@undefined\protect\let\protect\relax\fi
158 \ifx\documentclass\@undefined
159   \def\xstring{\string\string\string}
160 \else
161   \let\xstring\string
162 \fi
163 \xdef\bbl@language@stack{}
164 \def\bbl@push@language{%
165   \xdef\bbl@language@stack{\language+\bbl@language@stack}%
166 }
167 \def\bbl@pop@lang#1+#2-#3{%
168   \def\language{#1}\xdef#3{#2}%
169 }
170 \def\bbl@pop@language{%
171   \expandafter\bbl@pop@lang\bbl@language@stack-\bbl@language@stack
172   \expandafter\bbl@set@language\expandafter{\language}%
173 }
174 \expandafter\def\csname selectlanguage \endcsname#1{%
175   \bbl@push@language
176   \aftergroup\bbl@pop@language
177   \bbl@set@language{#1}}
178 \def\bbl@set@language#1{%
179   \edef\language{%
180     \ifnum\escapechar=\expandafter'\string#1\@empty
181       \else \string#1\@empty\fi}%
182   \select@language{\language}%
183   \if@filesw
184     \protected@write\@auxout{}\{\string\select@language{\language}}}%
185     \addtocontents{toc}{\xstring\select@language{\language}}}%
186     \addtocontents{lof}{\xstring\select@language{\language}}}%
187     \addtocontents{lot}{\xstring\select@language{\language}}}%
188   \fi}
189 \def\select@language#1{%
190   \expandafter\ifx\csname l@#1\endcsname\relax
191     \@nolanerr{#1}%
192   \else
193     \expandafter\ifx\csname date#1\endcsname\relax
194       \@noopterr{#1}%
195     \else
196       \bbl@patterns{\language}%
197       \originalTeX
198       \expandafter\def\expandafter\originalTeX
199         \expandafter{\csname noextras#1\endcsname

```

```

200             \let\originalTeX\@empty}%
201     \languageshorthands{none}%
202     \babel@beginsave
203     \csname captions#1\endcsname
204     \csname date#1\endcsname
205     \csname extras#1\endcsname\relax
206     \babel@savevariable\lefthyphenmin
207     \babel@savevariable\righthyphenmin
208     \expandafter\ifx\csname #1hyphenmins\endcsname\relax
209         \set@hyphenmins\tw@\thr@@\relax
210     \else
211         \expandafter\expandafter\expandafter\set@hyphenmins
212         \csname #1hyphenmins\endcsname\relax
213     \fi
214     \fi
215 \fi}
216 \long\def\otherlanguage#1{%
217     \csname selectlanguage \endcsname{#1}%
218     \ignorespaces
219 }
220 \long\def\endotherlanguage{%
221     \originalTeX
222     \global\@ignoretrue\ignorespaces
223 }
224 \expandafter\def\csname otherlanguage*\endcsname#1{%
225     \foreign@language{#1}%
226 }
227 \expandafter\def\csname endotherlanguage*\endcsname{%
228     \csname noextras\language\endcsname
229 }
230 \def\foreignlanguage{\protect\csname foreignlanguage \endcsname}
231 \expandafter\def\csname foreignlanguage \endcsname#1#2{%
232     \begingroup
233         \originalTeX
234         \foreign@language{#1}%
235         #2%
236         \csname noextras#1\endcsname
237     \endgroup
238 }
239 \def\foreign@language#1{%
240     \def\language{#1}%
241     \expandafter\ifx\csname l@#1\endcsname\relax
242         \@nolanerr{#1}%
243     \else
244         \bbl@patterns{\language}%
245         \languageshorthands{none}%
246         \csname extras#1\endcsname
247         \expandafter\ifx\csname #1hyphenmins\endcsname\relax
248             \set@hyphenmins\tw@\thr@@\relax
249         \else

```

```

250     \expandafter\expandafter\expandafter\set@hyphenmins
251     \csname #1hyphenmins\endcsname\relax
252   \fi
253 \fi
254 }
255 \def\bbl@patterns#1{%
256   \language=\expandafter\ifx\csname l@#1:\f@encoding\endcsname\relax
257     \csname l@#1\endcsname
258   \else
259     \csname l@#1:\f@encoding\endcsname
260   \fi\relax

```

With Lua_T_EX, 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 Lua_T_EX.

```

261   \ifx\directlua@\undefined\else
262     \ifx\directlua\relax\else
263       \ifcsname lu@texhyphen@loaded@the\language\endcsname \else
264         \global\@namedef{lu@texhyphen@loaded@the\language}{}%
265         \directlua{
266           if not luatexhyphen then
267             dofile(assert(kpse.find_file("luatex-hyphen.lua")))
268           end
269           luatexhyphen.loadlanguage("\luatexluaescapestring{#1}",
270             \the\language)}%
271       \fi
272     \fi
273   \fi
274 }
275 \def\hyphenrules#1{%
276   \expandafter\ifx\csname l@#1\endcsname@\undefined
277     \@nolanerr{#1}%
278   \else
279     \bbl@patterns{#1}%
280     \languageshorthands{none}%
281     \expandafter\ifx\csname #1hyphenmins\endcsname\relax
282       \set@hyphenmins\tw@\thr@@\relax
283     \else
284       \expandafter\expandafter\expandafter\set@hyphenmins
285       \csname #1hyphenmins\endcsname\relax
286     \fi
287   \fi
288 }
289 \def\endhyphenrules{}
290 \def\providehyphenmins#1#2{%
291   \expandafter\ifx\csname #1hyphenmins\endcsname\relax
292     \@namedef{#1hyphenmins}{#2}%
293   \fi}

```



```

294 \def\set@hyphenmins#1#2{\lefthyphenmin#1\righthyphenmin#2}
295 \def\LdfInit{%
296   \chardef\atcatcode=\catcode'\@
297   \catcode'\@=11\relax
298   \input babel.def\relax
299   \catcode'\@=\atcatcode \let\atcatcode\relax
300   \LdfInit}
301 \ifx\originalTeX\@undefined\let\originalTeX\@empty\fi
302 \ifx\babel@beginsave\@undefined\let\babel@beginsave\relax\fi
303 \ifx\PackageError\@undefined
304   \def\@nolanerr#1{%
305     \errhelp{Your command will be ignored, type <return> to proceed}%
306     \errmessage{You haven't defined the language #1\space yet}}
307   \def\@nopatterns#1{%
308     \message{No hyphenation patterns were loaded for}%
309     \message{the language '#1'}%
310     \message{I will use the patterns loaded for \string\language=0
311       instead}}
312   \def\@noopterr#1{%
313     \errmessage{The option #1 was not specified in \string\usepackage}
314     \errhelp{You may continue, but expect unexpected results}}
315   \def\@activated#1{%
316     \wlog{Package babel Info: Making #1 an active character}}
317 \else
318   \newcommand*{\@nolanerr}[1]{%
319     \PackageError{babel}%
320       {You haven't defined the language #1\space yet}%
321       {Your command will be ignored, type <return> to proceed}}
322   \newcommand*{\@nopatterns}[1]{%
323     \PackageWarningNoLine{babel}%
324       {No hyphenation patterns were loaded for\MessageBreak
325         the language '#1'\MessageBreak
326         I will use the patterns loaded for \string\language=0
327         instead}}
328   \newcommand*{\@noopterr}[1]{%
329     \PackageError{babel}%
330       {You haven't loaded the option #1\space yet}%
331       {You may proceed, but expect unexpected results}}
332   \newcommand*{\@activated}[1]{%
333     \PackageInfo{babel}{%
334       Making #1 an active character}}
335 \fi
336 \def\process@line#1#2 #3/{%
337   \ifx=#1
338     \process@synonym#2 /
339   \else
340     \process@language#1#2 #3/%
341   \fi
342 }
343 \toks@{}

```

```

344 \def\process@synonym#1 /{%
345   \ifnum\last@language=\m@ne
346     \expandafter\chardef\csname l@#1\endcsname0\relax
347     \wlog{\string\l@#1=\string\language0}
348     \toks@\expandafter{\the\toks@
349       \expandafter\let\csname #1hyphenmins\expandafter\endcsname
350       \csname\language\hyphenmins\endcsname}%
351   \else
352     \expandafter\chardef\csname l@#1\endcsname\last@language
353     \wlog{\string\l@#1=\string\language\the\last@language}
354     \expandafter\let\csname #1hyphenmins\expandafter\endcsname
355     \csname\language\hyphenmins\endcsname
356   \fi
357 }
358 \def\process@language#1 #2 #3/{%
359   \expandafter\addlanguage\csname l@#1\endcsname
360   \expandafter\language\csname l@#1\endcsname
361   \def\language{#1}%

```

In the LuaTeX case, we have to decide whether to load the language now.

```

362 \ifx\directlua@\undefined
363   \global\toks8\expandafter{\the\toks8#1, }%
364 \else
365   \directlua{
366     if not luatexhyphen then
367       dofile(assert(kpse.find_file("luatex-hyphen.lua")))
368     end
369     processnow = (tex.language == 0) or
370       (luatexhyphen.lookupname("\luatexluaescapestring{#1}") == nil)}%
371   \ifnum0=\directlua{tex.sprint(processnow and "0" or "1")}\relax
372   \global\toks8\expandafter{\the\toks8#1, }%
373   \global\@namedef{lu@texhyphen@loaded@the\language}{}%
374   \fi
375 \fi
376 \begingroup
377   \bbl@get@enc#1:\@@@
378   \ifx\bbl@hyph@enc@\empty
379   \else
380     \fontencoding{\bbl@hyph@enc}\selectfont
381   \fi
382   \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.

```

383   \ifx\directlua@\undefined
384     \input #2\relax
385   \else

```

```

386     \ifnum0=\directlua{tex.sprint(processnow and "0" or "1")}\relax
387     \input #2\relax
388     \fi
389 \fi
390 \ifnum\lefthyphenmin=\m@ne
391 \else
392     \expandafter\xdef\csname #1hyphenmins\endcsname{%
393         \the\lefthyphenmin\the\righthyphenmin}%
394 \fi
395 \endgroup
396 \ifnum\the\language=\z@
397     \expandafter\ifx\csname #1hyphenmins\endcsname\relax
398         \set@hyphenmins\tw@\thr@@\relax
399     \else
400         \expandafter\expandafter\expandafter\set@hyphenmins
401         \csname #1hyphenmins\endcsname
402     \fi
403     \the\toks@
404 \fi
405 \toks@{}%
406 \def\bbl@tempa{#3}%
407 \ifx\bbl@tempa\@empty
408 \else
409     \ifx\bbl@tempa\space
410 \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.

```

411     \ifx\directlua\@undefined
412         \input #3\relax
413     \else
414         \ifnum0=\directlua{tex.sprint(processnow and "0" or "1")}\relax
415         \input #3\relax
416     \fi
417     \directlua{processnow = nil}%
418 \fi
419 \fi
420 \fi
421 }
422 \def\bbl@get@enc#1:#2\@@@{%
423 \def\bbl@tempa{#1}%
424 \def\bbl@tempb{#2}%
425 \ifx\bbl@tempb\@empty
426     \let\bbl@hyph@enc\@empty
427 \else
428     \bbl@get@enc#2\@@@
429     \edef\bbl@hyph@enc{\bbl@tempa}%

```

```

430 \fi}
431 \openin1 = language.dat
432 \ifeof1
433 \message{I couldn't find the file language.dat,\space
434         I will try the file hyphen.tex}
435 \input hyphen.tex\relax
436 \else
437 \last@language@m@ne
438 \loop
439 \endlinechar@m@ne
440 \read1 to \bbl@line
441 \endlinechar'\^^M
442 \ifx\bbl@line\@empty
443 \else
444 \edef\bbl@line{\bbl@line\space/}%
445 \expandafter\process@line\bbl@line
446 \fi
447 \iftrue \csname fi\endcsname
448 \csname if\ifeof1 false\else true\fi\endcsname
449 \repeat
450 \language=0
451 \fi
452 \closein1
453 \let\process@language\@undefined
454 \let\process@synonym\@undefined
455 \let\process@line\@undefined
456 \let\bbl@tempa\@undefined
457 \let\bbl@tempb\@undefined
458 \let\bbl@eqq\@undefined
459 \let\bbl@line\@undefined
460 \let\bbl@get@enc\@undefined
461 \ifx\addto@hook\@undefined
462 \else
463 \expandafter\addto@hook\expandafter\everyjob\expandafter{%
464 \expandafter\typeout\expandafter{\the\toks8 loaded.}}
465 \fi

```