

# The `hyph-utf8` package and hyphenation with $\TeX$

Maintainers of the `hyph-utf8` package and patterns collection:

- *Mojca Miklavc, Arthur Reutenauer* (core, patterns)
- *Manuel Pégourié-Gonnard, Khaled Hosny* (Lua $\TeX$  support)
- *Élie Roux* (no longer active)

Date:

2017-05-18

Abstract:

*In 2008 all the existing hyphenation patterns from  $\TeX$  distributions have been collected in a single package `hyph-utf8`, converted into UTF-8 encoding and adapted for use in different  $\TeX$  engines. The patterns can be used directly by Unicode-aware engines such as Lua $\TeX$  and Xe $\LaTeX$ , and there is a mechanism to convert the patterns to the appropriate 8-bit encoding when used with p $\TeX$ , pdf $\TeX$  or Knuth's  $\TeX$ .*

Table of Contents:

<b>1</b>	<b>USING HYPHENATION PATTERNS</b> .....	<b>2</b>
1.1	Plain $\TeX$ .....	2
1.2	L <sup>A</sup> $\TeX$ .....	2
	<i>L<sup>A</sup><math>\TeX</math> with Babel</i> .....	2
	<i>L<sup>A</sup><math>\TeX</math> with Polyglossia</i> .....	2
	<i>Low-level commands</i> .....	2
1.3	Con $\TeX$ t .....	3
	<i>Con<math>\TeX</math>t MkII</i> .....	3
1.4	Some advanced examples .....	4
	<i>Example for Polyglossia</i> .....	4
<b>2</b>	<b>LIST OF SUPPORTED LANGUAGES</b> .....	<b>5</b>

# 1 Using hyphenation patterns

## 1.1 Plain T<sub>E</sub>X

In engines that support  $\epsilon$ -T<sub>E</sub>X you can select the desired hyphenation patterns with:

```
\uselanguage{langname}
```

where `langname` is the string identifying a particular hyphenation file in `language.def` (see Section 2).

## 1.2 L<sup>A</sup>T<sub>E</sub>X

### 1.2.1 L<sup>A</sup>T<sub>E</sub>X with Babel

You can switch the language in L<sup>A</sup>T<sub>E</sub>X with:

```
\usepackage[languagename]{babel}
```

In 8-bit engines you also need to make sure that you load the proper font encoding which supports all the characters used in the language of your choice, for example:

```
\usepackage[T1]{fontenc}
```

*N.B.:* You can use Babel with any T<sub>E</sub>X engine, however it has never been properly adapted to work well with Unicode engines. If you are using X<sub>Y</sub>T<sub>E</sub>X it is advisable to use Polyglossia instead.

### 1.2.2 L<sup>A</sup>T<sub>E</sub>X with Polyglossia

Polyglossia should be the preferred choice when using XeL<sup>A</sup>T<sub>E</sub>X. It doesn't support LuaL<sup>A</sup>T<sub>E</sub>X yet, but it is planned to extend it in future.

```
\usepackage{polyglossia}
\setmainlanguage[optional settings]{langname}
\setotherlanguages{otherlangname}

\begin[optional settings]{otherlangname} ... \end{otherlangname}
```

See Polyglossia manual for extensive list of options.

### 1.2.3 Low-level commands

Since Babel's `hyphen.cfg` is built into the L<sup>A</sup>T<sub>E</sub>X format, hyphenation patterns can be used without even loading Babel or Polyglossia. At the low-level this simply corresponds to defining

```
\language=\l@<langname>
```

The user command is supposed to be

```
\hyphenrules{langname}
```

or

```
\begin{hyphenrules}{langname} ... \end{hyphenrules}.
```

and *should* work with any flavour of L<sup>A</sup>T<sub>E</sub>X, however we couldn't make it work.

## 1.3 ConT<sub>E</sub>Xt

ConT<sub>E</sub>Xt doesn't load patterns for all the language that `hyph-utf8` provides. If you miss any language, please contact the mailing list. The general syntax for supported languages is the following:

```
% language of the main document
\mainlanguage[language]

% to switch to another language locally
{\language[otherlanguage] language of some short fragment}
```

You can use full language name or the two-letter language code.

### 1.3.1 ConT<sub>E</sub>Xt MkII

When using ConT<sub>E</sub>Xt MkII the EC/T1 font encoding is used by default already, but you might need to change the encoding when using Polish, languages written in Cyrillic scripts, etc. For example:

```
\usetypescript[iwona][qx]
\setupbodyfont[iwona]
\mainlanguage[polish]
```

ConT<sub>E</sub>Xt loads hyphenation patterns in several encodings. The Czech or Slovak patterns can be used with both EC and IL2 font encoding for example. The right hyphenation patterns will be chosen based on current font encoding.

## 1.4 Some advanced examples

### 1.4.1 Example for Polyglossia

```
\usepackage{polyglossia}
% the language used for main document
\setmainlanguage{asturian}
% American English with extended hyphenation patterns
\setotherlanguage[variant=usmax]{english}
% German with experimental patterns "ngerman-x-latest"
\setotherlanguage[spelling=new,latesthyphen=true]{german}
\setotherlanguages{spanish,catalan,french}

\begin{document}

Long Asturian text ... (Hyphenation for Asturian is not available,
but polyglossia automatically falls back on Catalan for now,
which seems to be a reasonable choice.)

\begin{german}
Deutscher Text ... (with the hyphenation patterns selected above:
"ngerman-x-latest")
\end{german}

\begin[script=fraktur,spelling=old]{german}
Deutlicher Text ... (set in Fraktur, with traditional hyphenation).
\end{german}

\end{document}
```

## 2 List of supported languages

<b>English</b>					
-	<b>english</b>	usenglish, USenglish, american			
en-us	usenglishmax				
en-gb	ukenglish	british, UKenglish			
<b>Afrikaans</b>					
af	afrikaans				
<b>Ancientgreek</b>					
grc	ancientgreek				
grc-x-ibycus	ibycus				
<b>Arabic</b>					
ar	arabic				
<b>Armenian</b>					
hy	armenian				
<b>Assamese</b>					
as	assamese				
<b>Basque</b>					
eu	basque				
<b>Belarusian</b>					
be	belarusian				
<b>Bengali</b>					
bn	bengali				
<b>Bulgarian</b>					
bg	bulgarian				
<b>Catalan</b>					
ca	catalan				
<b>Chinese</b>					
zh-latn-pinyin	pinyin				
<b>Church Slavonic</b>					
cu	churchslavonic				
<b>Coptic</b>					
cop	coptic				
<b>Croatian</b>					
hr	croatian				
<b>Czech</b>					
cs	czech				
<b>Danish</b>					
da	danish				
<b>Dutch</b>					
nl	dutch				
<b>Esperanto</b>					
eo	esperanto				
<b>Estonian</b>					
et	estonian				
<b>Ethiopic</b>					
mul-ethi	ethiopic	amharic, geez			
<b>Farsi</b>					
fa	farsi	persian			
<b>Finnish</b>					
fi	finnish				
<b>French</b>					
fr	french	patois, francais			
<b>Friulan</b>					
fur	friulan				
<b>Galician</b>					
gl	galician				
<b>Georgian</b>					
ka	georgian				
<b>German</b>					
de-1901	german				
de-1996	ngerman				
de-ch-1901	swissgerman				
<b>Greek</b>					
el-monoton	monogreek				
el-polyton	greek	polygreek			
<b>Gujarati</b>					
gu	gujarati				
<b>Hindi</b>					
hi	hindi				
<b>Hungarian</b>					
hu	hungarian				
<b>Icelandic</b>					
is	icelandic				
<b>Indonesian</b>					
id	indonesian				
<b>Interlingua</b>					
ia	interlingua				
<b>Irish</b>					
ga	irish				
<b>Italian</b>					
it	italian				
<b>Kannada</b>					
kn	kannada				
<b>Kurmanji</b>					
kmr	kurmanji				
<b>Latin</b>					
la	latin				
la-x-classic	classicalatin				
la-x-liturgic	liturgicallatin				
<b>Latvian</b>					
lv	latvian				
<b>Lithuanian</b>					
lt	lithuanian				
<b>Malayalam</b>					
ml	malayalam				

<b>Marathi</b>			<b>Serbian</b>		
mr	marathi		sr-latn	serbian	
<b>Mongolian</b>			sr-cyrl	serbianc	
mn-cyrl	mongolian		<b>Slovak</b>		
mn-cyrl-x-lmc	mongolianlmc		sk	slovak	
<b>Norwegian</b>			<b>Slovenian</b>		
nb	bokmal	norwegian, norsk	sl	slovenian	slovene
nn	nynorsk		<b>Spanish</b>		
<b>Occitan</b>			es	spanish	espanol
oc	occitan		<b>Swedish</b>		
<b>Oriya</b>			sv	swedish	
or	oriya		<b>Tamil</b>		
<b>Panjabi</b>			ta	tamil	
pa	panjabi		<b>Telugu</b>		
<b>Polish</b>			te	telugu	
pl	polish		<b>Thai</b>		
<b>Piedmontese</b>			th	thai	
pms	piedmontese		<b>Turkish</b>		
<b>Portuguese</b>			tr	turkish	
pt	portuguese	portuges	<b>Turkmen</b>		
<b>Romanian</b>			tk	turkmen	
ro	romanian		<b>Ukrainian</b>		
<b>Romansh</b>			uk	ukrainian	
rm	romansh		<b>Uppersorbian</b>		
<b>Russian</b>			hsb	uppersorbian	
ru	russian		<b>Welsh</b>		
<b>Sanskrit</b>			cy	welsh	
sa	sanskrit				

Babel defines a few more synonyms (which consequently only work in L<sup>A</sup>T<sub>E</sub>X):

<b>english</b>	canadian
<b>british</b>	australian, newzealand
<b>german</b>	austrian
<b>ngerman</b>	naustrian
<b>portuguese</b>	brazilian, brazil