

# hyph-utf8

Mojca Miklavc (current maintainer)  
Arthur Reutenauer (no longer active on hyph-utf8)

April 9, 2012

## Abstract

The hyph-utf8 package gathers all the existing hyphenation patterns for T<sub>E</sub>X, converted to UTF-8. They can be used directly by UTF-8-aware T<sub>E</sub>X engines such as LuaT<sub>E</sub>X and X<sub>Y</sub>T<sub>E</sub>X, and there is a mechanism to convert the patterns to some 8-bit encoding when used with pdfT<sub>E</sub>X or Knuth's T<sub>E</sub>X.

## List of supported languages

<b>English</b>			
-	english	usenglish, USenglish, american	
en-us	usenglishmax		
en-gb	ukenglish	british, UKenglish	
<b>Afrikaans</b>		<b>Esperanto</b>	
af	afrikaans	eo	esperanto
<b>Ancientgreek</b>		<b>Estonian</b>	
grc	ancientgreek	et	estonian
grc-x-ibycus	ibycus	<b>Ethiopic</b>	
<b>Arabic</b>		mul-ethi	ethiopic amharic, geez
ar	arabic	<b>Farsi</b>	
<b>Armenian</b>		fa	farsi persian
hy	armenian	<b>Finnish</b>	
<b>Assamese</b>		fi	finnish
as	assamese	<b>French</b>	
<b>Basque</b>		fr	french patois, francais
eu	basque	<b>Friulan</b>	
<b>Bengali</b>		fur	friulan
bn	bengali	<b>Galician</b>	
<b>Bulgarian</b>		gl	galician
bg	bulgarian	<b>German</b>	
<b>Catalan</b>		de-1901	german
ca	catalan	de-1996	ngerman
<b>Chinese</b>		de-ch-1901	swissgerman
zh-latn-pinyin	pinyin	<b>Greek</b>	
<b>Coptic</b>		el-monoton	monogreek
cop	coptic	el-polyton	greek polygreek
<b>Croatian</b>		<b>Gujarati</b>	
hr	croatian	gu	gujarati
<b>Czech</b>		<b>Hindi</b>	
cs	czech	hi	hindi
<b>Danish</b>		<b>Hungarian</b>	
da	danish	hu	hungarian
<b>Dutch</b>		<b>Icelandic</b>	
nl	dutch	is	icelandic

<b>Indonesian</b>			<b>Portuguese</b>		
id	indonesian		pt	portuguese	portuges
<b>Interlingua</b>			<b>Romanian</b>		
ia	interlingua		ro	romanian	
<b>Irish</b>			<b>Romansh</b>		
ga	irish		rm	romansh	
<b>Italian</b>			<b>Russian</b>		
it	italian		ru	russian	
<b>Kannada</b>			<b>Sanskrit</b>		
kn	kannada		sa	sanskrit	
<b>Kurmanji</b>			<b>Serbian</b>		
kmr	kurmanji		sr-latn	serbian	
<b>Lao</b>			sr-cyrl	serbianc	
lo	lao		<b>Slovak</b>		
<b>Latin</b>			sk	slovak	
la	latin		<b>Slovenian</b>		
<b>Latvian</b>			sl	slovenian	slovene
lv	latvian		<b>Spanish</b>		
<b>Lithuanian</b>			es	spanish	espanol
lt	lithuanian		<b>Swedish</b>		
<b>Malayalam</b>			sv	swedish	
ml	malayalam		<b>Tamil</b>		
<b>Marathi</b>			ta	tamil	
mr	marathi		<b>Telugu</b>		
<b>Mongolian</b>			te	telugu	
mn-cyrl	mongolian		<b>Turkish</b>		
mn-cyrl-x-lmc	mongolianlmc		tr	turkish	
<b>Norwegian</b>			<b>Turkmen</b>		
nb	bokmal	norwegian, norsk	tk	turkmen	
nn	nynorsk		<b>Ukrainian</b>		
<b>Oriya</b>			uk	ukrainian	
or	oriya		<b>Uppersorbian</b>		
<b>Panjabi</b>			hsb	uppersorbian	
pa	panjabi		<b>Welsh</b>		
<b>Polish</b>			cy	welsh	
pl	polish				

# Using hyphenation patterns

## 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.dat` and can be taken from table on the first two pages.

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

Since Babel's `hyphen.cfg` is built in the XeL<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.

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

You can use Babel with any T<sub>E</sub>X engine, however it is currently unmaintained and has never been adapted to work well with Unicode engines. If you are using XeT<sub>E</sub>X please use Polyglossia instead.

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

## 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.

## 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]

{\language[otherlanguage] language of some short fragment}
```

You can use full language name or language code. When using ConT<sub>E</sub>Xt MKII you might need to select the appropriate font encoding for Cyrillic scripts, Polish and some other languages:

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

ConT<sub>E</sub>Xt loads hyphenation patterns in several encodings, so that you can for example use Czech patterns with either ec or il2 font encodings. The right hyphenation patterns will be chosen based on current font encoding.

## More examples

### 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}
Deutfcher Text ... (set in Fraktur, with traditional hyphenation).
\end{german}
```

```
\end{document}
```