

# The hyph-utf8 package and hyphenation with T<sub>E</sub>X

Maintainers of the hyph-utf8 package and collectors of patterns:

- Mojca Miklavec, Arthur Reutenauer
- With contributions by Khaled Hosny, Manuel Pégourié-Gonnard, Élie Roux

Main description:

June 2011

Latest editorial change:

16 Mars 2018

Abstract:

In 2008 all the existing hyphenation patterns from T<sub>E</sub>X distributions have been collected in a single package *hyph-utf8*, converted into UTF-8 encoding and adapted for use in different T<sub>E</sub>X engines. The patterns can be used directly by Unicode-aware engines such as *LuaT<sub>E</sub>X* and *X<sub>E</sub>T<sub>E</sub>X*, and there is a mechanism to convert the patterns to the appropriate 8-bit encoding when used with *pT<sub>E</sub>X*, *pdfT<sub>E</sub>X* or Knuth's T<sub>E</sub>X.

Table of Contents:

<b>1</b>	<b>USING HYPHENATION PATTERNS</b>	2
1.1	Plain T <sub>E</sub> X	2
1.2	L <sup>A</sup> T <sub>E</sub> X	2
	L <sup>A</sup> T <sub>E</sub> X with Babel	2
	L <sup>A</sup> T <sub>E</sub> X with Polyglossia	2
	Low-level commands	2
1.3	ConT <sub>E</sub> Xt	3
	ConT <sub>E</sub> Xt MkII	3
1.4	Some advanced examples	4
	Example for Polyglossia	4
<b>2</b>	<b>LIST OF SUPPORTED LANGUAGES</b>	5

# 1 Using hyphenation patterns

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

In engines that support ε-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>E</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 ConTeXt

ConTeXt 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 ConTeXt MkII

When using ConTeXt 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]
```

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

\end{document}
```

## 2 List of supported languages

For several languages, there is additional documentation in a separate file: see

- For German, `dehyph-exptl.pdf`
- For Spanish, `division.pdf`

### English

-	english	usenglish, USenglish, american
en-us	usenglishmax	
en-gb	ukenglish	british, UKenglish

### Afrikaans

af	afrikaans
----	-----------

### Ancientgreek

grc	ancientgreek
grc-x-ibycus	ibycus

### Arabic

ar	arabic
----	--------

### Armenian

hy	armenian
----	----------

### Assamese

as	assamese
----	----------

### Basque

eu	basque
----	--------

### Belarusian

be	belarusian
----	------------

### Bengali

bn	bengali
----	---------

### Bulgarian

bg	bulgarian
----	-----------

### Catalan

ca	catalan
----	---------

### Chinese

zh-latn-pinyin	pinyin
----------------	--------

### Church Slavonic

cu	churchslavonic
----	----------------

### Coptic

cop	coptic
-----	--------

### Croatian

hr	croatian
----	----------

### Czech

cs	czech
----	-------

### Danish

da	danish
----	--------

### Dutch

nl	dutch
----	-------

### Esperanto

eo	esperanto
----	-----------

### Estonian

et	estonian
----	----------

### Ethiopic

mul-ethi	ethiopic	amharic, gezz
----------	----------	---------------

### Farsi

fa	farsi	persian
----	-------	---------

### Finnish

fi	finnish
----	---------

<b>French</b>			<b>Marathi</b>	
fr	french	patois, francais	mr	marathi
<b>Friulan</b>			<b>Mongolian</b>	
fur	friulan		mn-cyrl	mongolian
<b>Galician</b>			mn-cyrl-x-lmc	mongolianlmc
gl	galician		<b>Norwegian</b>	
<b>Georgian</b>			nb	bokmal
ka	georgian		nn	nynorsk
<b>German</b>			<b>Occitan</b>	
de-1901	german		oc	occitan
de-1996	ngerman		<b>Oriya</b>	
de-ch-1901	swissgerman		or	oriya
<b>Greek</b>			<b>Punjabi</b>	
el-monoton	monogreek		pa	panjabi
el-polyton	greek	polygreek	<b>Polish</b>	
<b>Gujarati</b>			pl	polish
gu	gujarati		<b>Piedmontese</b>	
<b>Hindi</b>			pms	piedmontese
hi	hindi		<b>Portuguese</b>	
<b>Hungarian</b>			pt	portuguese
hu	hungarian		<b>Romanian</b>	
<b>Icelandic</b>			ro	romanian
is	icelandic		<b>Romansh</b>	
<b>Indonesian</b>			rm	romansh
id	indonesian		<b>Russian</b>	
<b>Interlingua</b>			ru	russian
ia	interlingua		<b>Sanskrit</b>	
<b>Irish</b>			sa	sanskrit
ga	irish		<b>Serbian</b>	
<b>Italian</b>			sr-latn	serbian
it	italian		sr-cyrl	serbanc
<b>Kannada</b>			<b>Slovak</b>	
kn	kannada		sk	slovak
<b>Kurmanji</b>			<b>Slovenian</b>	
kmr	kurmanji		sl	slovenian
<b>Latin</b>			<b>Spanish</b>	
la	latin		es	spanish
la-x-classic	classiclatin		<b>Swedish</b>	
la-x-liturgic	liturgicallatin		sv	swedish
<b>Latvian</b>			<b>Tamil</b>	
lv	latvian		ta	tamil
<b>Lithuanian</b>			<b>Telugu</b>	
lt	lithuanian		te	telugu
<b>Malayalam</b>			<b>Thai</b>	
ml	malayalam		th	thai
			<b>Turkish</b>	
			tr	turkish
			<b>Turkmen</b>	
			tk	turkmen
			<b>Ukrainian</b>	
			uk	ukrainian
			<b>Uppersorbian</b>	
			hsb	uppersorbian
			<b>Welsh</b>	
			cy	welsh

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, nswissgerman
<b>portuguese</b>	brazilian, brazil