

ucs.sty - Unicode Support

Dominique P. G. Unruh

2004/10/17

Contents

1	Usage	1
1.1	Special options	2
1.2	Normal options	3
1.3	Combining mode	4
1.4	Defining unicode data	5
1.5	Known problems	5
2	Thanks	6

1 Usage

Simply use `\usepackage{ucs}` and `\usepackage[utf8x]{inputenc}`, then you will be able to write your LaTeX-Documents in UTF-8.

You can access a Unicode character with `\unichar{<code>}`, even when the active input encoding is not `utf8`.

An unicode character can have an default glyph macro and several glyph macros associated with options. If one of these options is set, the associated macro is used, otherwise the default macro. If several associated options are set, an error is yielded. You may activate an option *<name>* by including it in the option list while loading the ucs package, or by using `\SetUnicodeOption{<name>}`. To deactivate an option, prefix its name by `no`. Note that you must load `ucs.sty` before `\usepackage[utf8x]{inputenc}` if you want to supply options. Any option which you want to use must be used at least once in the preamble.

`\SetUnicodeOption`

When you activate an option, you can supply a priority as optional argument. If there are several glyphs for a given code position, the one having the option with the highest priority is taken (an error is yielded in case of ambiguity). If you do

not supply a priority, 100 is taken as default. “Normal” glyphs are associated with the option `default`, which is initially activated with a priority of 0.

NB: UTF-8 characters are interpreted by \TeX as a sequence of commands, so don’t use calls like `\macro ä` instead of `\macro{ä}` (this does not apply to ASCII characters).

This input encoding does not change the fontencoding automatically. You can use the package `autofo` for that purpose (<http://www.unruh.de/DniQ/latex/unicode/ucs/contrib/autofo.sty>).

1.1 Special options

Several options have a special hardcoded meaning:

<code>combine</code>	<ul style="list-style-type: none"> • <code>combine</code>: Activates combining mode. See section 1.3.
<code>default</code>	<ul style="list-style-type: none"> • <code>default</code>: This option contains all characters, which are not explicitly associated with another option. But see also the option <code>document</code>. This option is activated per default, but has priority 0, i.e. any other activated option is preferred, unless its priority is explicitly given to be smaller.
<code>document</code>	<ul style="list-style-type: none"> • <code>document</code>: Every character you define in your document using <code>\DeclareUnicodeCharacter</code> has the option <code>document</code>. This option is activated per default having the priority 1000, therefore manually declared characters take precedence over all other characters, unless some other option has explicitly gotten a higher priority.
<code>fasterrors</code>	<ul style="list-style-type: none"> • <code>fasterrors</code>: When used, the name of an unicode character is not included in error messages any more, which runs much faster.
<code>graphics</code>	<ul style="list-style-type: none"> • <code>graphics</code>: When used, unknown characters are replaced by GIFs downloaded from unicode.org. Commands to download and convert these are executed if <code>-shell-escape</code> is passed to \LaTeX, otherwise they are proposed to the user in a warning message. A UNIX-machine supporting the commands <code>wget</code>, <code>giftopnm</code> and <code>pnmtops</code> is assumed.
<code>savemem</code>	<ul style="list-style-type: none"> • <code>savemem</code>: When used, only the character needed at the moment is loaded, not a whole page. This slows down operation, but saves space in the \TeX-pool, especially with sparsely used character set like kanji. Use this, if you get an out of pool error or similar from \TeX. If you change the state of this option during the run, you may get unexpected results.
<code>warnunknown</code>	<ul style="list-style-type: none"> • <code>warnunknown</code>: When used, an unknown unicode character does not generate an error, but a warning.

1.2 Normal options

The options described here are—strictly spoken—not part of `ucs.sty`, but are defined by the unicode data files. They are included here for convenience.

<code>autogenerated</code>	<ul style="list-style-type: none"> • autogenerated: This enables the characters, which are autogenerated as composition of other characters according to the informations in the <code>UnicodeData.txt</code>. These may or may not look good. You may have to define the <code>\unicodetextcircle</code>, <code>\unicodetextsquare</code>, <code>\unicodetextvertical</code>, <code>\unicodetextwide</code> and <code>\unicodetextsmall</code> macros to let all autogenerated macros work. Furthermore you may have to set some other options, when the autogenerated characters are build out of characters, which are not in the default set.
<code>ckkgb5</code>	<ul style="list-style-type: none"> • ckkgb5: See the explanation for <code>ckkjis</code> below and substitute C40 by C00, JIS by BIG-5 and <code>kanji48</code> by some appropriate font which has BIG-5 encoding (e.g. one of the "Arphic AR PL * Big5" fonts).
<code>ckkgb</code>	<ul style="list-style-type: none"> • ckkgb: See the explanation for <code>ckkjis</code> below and substitute C40 by C10, JIS by GB and <code>kanji48</code> by some appropriate font which has GB encoding (e.g. one of the "Arphic AR PL * GB" fonts).
<code>ckkhangul</code>	<ul style="list-style-type: none"> • ckkhangul: See the explanation for <code>ckkjis</code> below and substitute C40 by C61, JIS by "KSC5601 hangul syllables" and <code>kanji48</code> by some appropriate font which has KSC5601 encoding and hangul syllables (e.g. the <code>han</code> or the <code>han1</code> font from CJK-\LaTeX).
<code>ckkjis</code>	<ul style="list-style-type: none"> • ckkjis: This enables the use of C40 (JIS) or C42 (JISdnp) encoded fonts. You need to have the <code>c40*.fd</code> files which are contributed with the package CJK and the <code>kanji48</code> font installed for this. Further you have to load the fontencoding C40 (an option to the package <code>fontenc</code>). <p>It is not necessary to load the package CJK. If you want to use it nevertheless take care of the following:</p> <ul style="list-style-type: none"> – Load CJK before <code>fontenc</code>, or quite strange errors will occur. – Load CJK with option <code>encapsulated</code>, or it will overwrite some of the UTF8 input encoding. – Don't use the CJK environment, it destroys the input encoding. Use <code>ucjk</code> instead, which is a patched version and takes no arguments.
<code>fullmathletters</code>	<ul style="list-style-type: none"> • fullmathletters: This option has been removed. Replace all occurrences by <code>mathletters</code>.
<code>mathletters</code>	<ul style="list-style-type: none"> • mathletters: When using this option is set, some unicode code characters like greek or some hebrew letters generate the math mode glyphs. This option is disabled by default, because using math greek in a normal text does not look good. But you may set it in <code>\everymath</code> and <code>\everydisplay</code> and such enable the use of unicode characters in math mode. You can also

use this to get a poor man's greek, it is however recommended to use the `cb` fonts instead.

- `postscript`
 - **postscript**: This option enables use of postscript, e.g. of postscript fonts. Some DVI-viewers may have problems with documents using this option, but most viewers can handle it correctly.
- `privatecsur`
 - **privatecsur**: This option enables use of characters in the private area according to the mapping by the ConScript Unicode Registry (<http://www.evertype.com/standards/csur/> and also <http://home.ccil.org/~cowan/csur/index.html>).
- `tipa`
 - **tipa**: This enables the use of the macros in the `tipa` package to display IPA symbols.

1.3 Combining mode

In some cases, Unicode documents contain sequences like U+0063 LATIN SMALL LETTER C U+0301 COMBINING ACUTE ACCENT (producing \acute{c}). In order to typeset them correctly, we cannot render U+0063 LATIN SMALL LETTER c immediately when it occurs, instead we have to wait whether some combining character is going to follow.

`combine` To enable this way of parsing in `ucs.sty`, we have to use the option `combine`. As long as this option is in effect, characters are not immediately rendered, but stored in a token register until they are output via

- resetting of the option (`\SetUnicodeOption{nocombine}`) or

`\unicodecombine`

- the command `\unicodecombine`.

Since characters handled by `ucs.sty` are not immediately output while ASCII characters are handled by \TeX and directly rendered, you should not mix ASCII characters and non-ASCII characters while this option is in effect. Thus to obtain the above glyph, you cannot use

`\SetUnicodeOption{combine}c\acute{\SetUnicodeOption{nocombine}}`

(where $\acute{\}$ is U+0301 COMBINING ACUTE ACCENT), instead you can use one of the following constructions:

- `\SetUnicodeOption{combine}\unichar{"63}\acute{\SetUnicodeOption{nocombine}}`

`\unicodevirtual`

- `\SetUnicodeOption{combine}\unicodevirtual{c}\acute{\SetUnicodeOption{nocombine}}` (`\unicodevirtual` takes arbitrary \LaTeX code and inserts it, as though it was a Unicode character; do not use Unicode characters inside `\unicodevirtual`).

- `\SetUnicodeOption{combine}\myverbatim|có|%`
`\SetUnicodeOption{nocombine}`
 where `\myverbatim`¹ is a command similar to `\verb`, but setting the catcodes of the ASCII characters to 13 (active) and then defining character no. *n* to expand to `\unichar{n}`.

In cases where you only want to render occasional words containing combining characters and *no* ASCII, you can use a macro like

```
\newcommand\combword[1]{\SetUnicodeOption{combine}#1%
\SetUnicodeOption{nocombine}}
```

and then simply typeset the concerning word as an argument to `\combword`.

1.4 Defining unicode data

A unicode character may be defined by

```
\DeclareUnicodeCharacter{<code>}{<macro>}
```

or, when it is to be associated with a special option, by using

```
\DeclareUnicodeCharacterAsOptional{<code>}{<option>}{<macro>}
```

where `<code>` is the unicode character number, `<option>` the associated option and `<macro>` the glyph's macro.

This definition is local.

In the automatically loaded data files `\uc@dc1c` should be used instead.

An option `<option>` can be defined by

```
\DeclareUnicodeOption[<pkg>]{<option>}.
```

If `<pkg>` is supplied, the option is set, if the package `<pkg>` is loaded.

You can add further packages, which automatically set an option, by

```
\LinkUnicodeOptionToPkg{<option>}{<pkg>}.
```

If a character *c* is unknown, it is looked up in the `uni-n.def`-file, whereby $n = \lfloor \frac{c}{256} \rfloor$. So characters which are not document specific, should be defined in those file. For generating them, you should use the program `makeunidef.pl`.

You may find `\dirtyunicode` and `\UnicodeNeeds` interesting for writing glyphs macros, see in the implementation section.

1.5 Known problems

Note, that if a character from some not yet loaded Unicode page appears, a file has to be loaded. If this appears inside a word, kerning and ligatures do not work at that position.

Further there are some commands, which expand their arguments in non-executing contexts, this makes it impossible for `ucs.sty` to load the character definition file

¹This command is not provided by `ucs.sty`, but some similar command may appear in future.

at this place if this has not yet been done. In this case the concerning character is replaced by some message that you have to use `\PrerenderUnicode{...}`.

Both problems can be solved by preloading the offending characters. If you have for example `U+03B1` GREEK SMALL LETTER ALPHA, which should be preloaded, simply use

<code>\PreloadUnicodePage</code>	<ul style="list-style-type: none"> • <code>\PreloadUnicodePage{3}</code> (the argument to <code>\PreloadUnicodePage</code> is $\lfloor \frac{n}{256} \rfloor$, where n is the number of the character (here <code>0x3B1</code> = 945)) or
<code>\PrerenderUnicode</code>	<ul style="list-style-type: none"> • <code>\PrerenderUnicode{α}</code> (the argument to <code>\PrerenderUnicode</code> can contain any L^AT_EX code, which is then rendered in an hbox, all still unknown characters are loaded and the result is thrown away; do not use e.g. <code>\footnote</code> or other commands which might not like to be executed several times).

2 Thanks

Thanks to...

- Michel Goossens who supplied many characters (e.g. Vietnamese, polytonic Greek),
- Manuel Kauers for testing my package before the first upload,
- Werner Lemberg, who wrote the CJK package, where I got the font definitions in `cencmn.tex` from,
- Karsten Tinnfeld for many of the glyph macros in `cyrillic.ucf`,
- Pablo Rodriguez for reporting many bugs, especially concerning the interaction between `ucs.sty` and other packages.
- Stefan Röhrich for testing my package before the first upload,
- the authors of all those many L^AT_EX-packages for different scripts.