

The `newunicodechar` package*

Enrico Gregorio
Enrico dot Gregorio at univr dot it

November 12, 2012

1 Introduction

When using Unicode input with \LaTeX it's not so uncommon to get an incomprehensible error message such as

```
Unicode char \u8:xxx not set up for use with LaTeX
```

where `xxx` may be the actual character we input or a combination of strange characters. This happens because the `utf8` option given to `inputenc` or `inputenx` defines the \LaTeX meaning of many Unicode characters, but, of course, not all of them.

For example, one might want to write some Latin words with prosodic marks, i.e., the diacritics that tell whether a vowel is long or short, in order to distinguish between ‘pōpūlus’ (people) and ‘pōpūlus’ (poplar), but using the actual Unicode characters that make the \LaTeX document easier to read; look at the following table, where on the left the input is via the \LaTeX Internal Character Representation (LICR) and on the right it's via Unicode characters,

LICR	Unicode
<code>p\u{o}p\u{u}lus</code>	pōpūlus
<code>p\={o}p\u{u}lus</code>	pōpūlus

and judge by yourselves which one is better. Unfortunately, the `utf8` option to `inputenc` doesn't define a meaning for `ō`, `ō`, and `ū`. As a matter of fact, only `ā` and `Ā` are defined, as they are used in the Romanian language.

One might resort to `\DeclareUnicodeCharacter` in the document's preamble, but this requires looking up at the (long) list of Unicode characters and jotting down the relevant numbers. For example, `ō` is `U+014F`, so the declaration

```
\DeclareUnicodeCharacter{014F}{\u{o}}
```

*This document corresponds to `newunicodechar` v1.1, dated 2012/11/12.

would do for δ .

The present package introduces a simpler interface that frees the user from the burden to look up in the tables: all it's needed is

```
\newunicodechar{\delta}{\u{o}}
```

You are not restricted to definition like this: for example, ♯ is Unicode U+266A, but you are not required to know it: if your editor can insert the character ♯, you may define its meaning by loading a package that provides it and say

```
\usepackage{wasysym}
\newunicodechar{\sharp}{\eighthnote}
```

A similar problem may arise even with X_YL^AT_EX. A frequently asked question on mailing lists or discussion groups is how to print some particular character in a different font than the main one of the document, say, for example, the Euro sign which, in some fonts, is horrible. The usual answer is to write something like

```
\newfontfamily{\eurofont}{\langle some font \rangle}
\catcode'\€=\active
\protected\def \€{{\eurofont\char'\€}}
```

which, for the average user, is somewhat scaring. With `newunicodechar` this may be simplified into

```
\newfontfamily{\eurofont}{\langle some font \rangle}
\newunicodechar{€}{{\eurofont\texteuro}}
```

2 Usage

The package requires the use of a Unicode engine, i.e., X_YL^AT_EX or Lua^AT_EX, or, with (pdf)L^AT_EX, the `inputenc` or `inputenx` package along with the `utf8` option. It won't work with the `utf8x` option that employs a completely different mechanism for parsing Unicode characters in (pdf)L^AT_EX. It should be said that `utf8x` defines many more characters than `utf8`, so that the present package wouldn't be needed.

Of course the L^AT_EX document must be written using a Unicode savvy editor.

verbose The package has only one option, `verbose`, which is off by default. If the package is called by `\usepackage[verbose]{newunicodechar}`, then the informative message on the log file will show the old definition along with the warning about the redefinition. Unfortunately this definition usually has a rather cryptic format; for example, redefining `ã` would print

```
Redefining Unicode character; it meant
*** \IeC {\u a} ***
before your redefinition on input line 22.
```

Call the package with this option if you are worried about what you are redefining, but the meaning of the Unicode character should correspond easily to just one LICR entry. This option does nothing when a Unicode engine is used, because no character is active initially (except for ~), so there should be non risk to redefine anything.

`\newunicodechar` The package provides only one command, `\newunicodechar`, which must be called with two arguments:

```
\newunicodechar{⟨char⟩}{⟨code⟩}
```

where `⟨char⟩` is the Unicode character to which we need to give a meaning and `⟨code⟩` is that meaning, that is the \LaTeX code that will be substituted to the character. Here is what's needed for the prosodic marks in Latin:

```
%\newunicodechar{Ā}{\u{A}} \newunicodechar{ā}{\u{a}}
\newunicodechar{Ē}{\u{E}} \newunicodechar{ē}{\u{e}}
\newunicodechar{Ī}{\u{I}} \newunicodechar{ī}{\u{i}}
\newunicodechar{Ō}{\u{O}} \newunicodechar{ō}{\u{o}}
\newunicodechar{Ū}{\u{U}} \newunicodechar{ū}{\u{u}}
\newunicodechar{Ā}{\={A}} \newunicodechar{ā}{\={a}}
\newunicodechar{Ē}{\={E}} \newunicodechar{ē}{\={e}}
\newunicodechar{Ī}{\={I}} \newunicodechar{ī}{\={i}}
\newunicodechar{Ō}{\={O}} \newunicodechar{ō}{\={o}}
\newunicodechar{Ū}{\={U}} \newunicodechar{ū}{\={u}}
```

The first line is commented out, because those characters are already defined. It doesn't hurt to give again a definition, \LaTeX will just warn about it.

Caution: when used with $X_{\text{pdf}}\LaTeX$ or $\text{Lua}\LaTeX$, there will be no warning, but the standard setup doesn't define any active character.

The first argument *must* consist of a single Unicode character, but the package checks for it and raises an error otherwise; it will raise an error also if the first argument consists of a plain ASCII character: defining them is not allowed in (pdf) \LaTeX and would break almost everything in $X_{\text{pdf}}\LaTeX$ or $\text{Lua}\LaTeX$.

3 An easier way?

One could dispense with this package, since the same effect may be obtained in the way we have already seen in $X_{\text{pdf}}\LaTeX$ or $\text{Lua}\LaTeX$; with (pdf) \LaTeX , calling `\newunicodechar{ū}{\={u}}` is equivalent to say

```
\makeatletter
\@namedef{u8:\detokenize{ū}}{\={u}}
\makeatother
```

4 Implementation

The usual presentation, that we repeat here for completeness.

```

\ProvidesFile{newunicodechar.dtx}
\NeedsTeXFormat{LaTeX2e}[2008/04/05]
\ProvidesPackage{newunicodechar}

```

The date for the format has been chosen to match the last version of the `utf8enc.dfu` file that provides definitions for Unicode characters.

Now the real macros. First of all we check that the typesetting engine is sufficiently recent to include ε -TeX extensions.

```

1 \ifundefined{eTeXversion}
2   {\PackageError{newunicodechar}{LaTeX engine too old, aborting}
3    {Please upgrade your TeX system}\@@end}{ }

```

4.1 Options

There's only one option.

```

4 \DeclareOption{verbose}{\let\nuc@verbose=T}
5 \ProcessOptions\relax

```

4.2 Error messages

```

6 \def\nuc@onebyteerr{\PackageError{newunicodechar}
7   {ASCII character requested}
8   {Only characters above U+007F may be defined; you asked
9    for\MessageBreak a plain ASCII character and your definition
10   has been ignored.}}
11 \def\nuc@emptyargerr{\PackageError{newunicodechar}
12   {Empty argument}
13   {You shouldn't write \protect\newunicodechar}{...}}
14 \def\nuc@invalidargerr{\PackageError{newunicodechar}
15   {Invalid argument}
16   {The first argument to \protect\newunicodechar\space is
17    either\MessageBreak too long or an invalid sequence of bytes}}

```

4.3 Unicode engines

In case we are running a Unicode engine (`xelatex` or `lualatex`), the definition of the main macro is easier: we check whether the character is above 127 and, in this case, we make it active expanding to the second argument. This definition of the main macro will be seen only if the engine has the `^^^` convention for inputting characters in hexadecimal format, so that `^^^0021` is one token and `\@gobble` will eat it up, making `\next` equal to `\@empty`; with an eight bit engine (`latex` or `pdflatex`), `\@gobble` will swallow only `^^^`.

In this case we define the main macro all in one swoop; first we check that the first argument is nonempty, then we act only if it consists of only one (character) token. Then if the charcode of this token is less than 127 we emit an error message; otherwise we activate the character and define its (protected) expansion to be the second argument. The last action, in this case, is to allow `\newunicodechar` only in the preamble. Then we do `\endinput`. If the engine is not Unicode savvy, everything up to the closing `\fi` is swallowed up.

`\newunicodechar` Here is the code for defining the Unicode engine version of the main macro. Notice that we try being on the safe side by not assuming any particular category code for `~`, because we restore it after giving the definition where we need it to be active.

```

18 \begingroup
19 \catcode'\^=7 \catcode30=12 \catcode'\!=12 % for safety
20 \edef\next{\@gobble^^^0021}
21 \expandafter\endgroup
22 \ifx\next\@empty % Start of code for Unicode engines
23 \chardef\nuc@atcode=\catcode'\~
24 \catcode'\~=\active
25 \def\newunicodechar#1#2{%
26   \if\relax\detokenize{#1}\relax
27     \nuc@emptyargerr
28   \else
29     \if\relax\detokenize\expandafter{\@cdr#1\@nil}\relax
30       \ifnum'#1>\string"7F
31         \catcode'#1=\active
32         \begingroup\lccode'\~='#1
33         \lowercase{\endgroup\protected\def~}{#2}%
34       \else
35         \nuc@onebyteerr
36       \fi
37     \else
38       \nuc@invalidargerr
39     \fi
40   \fi}
41 \catcode'\~=\nuc@atcode
42 \@onlypreamble\newunicodechar
43 \expandafter\endinput
44 \fi % End of code for Unicode engines

```

4.4 Eight bit engines

From now on we can assume an eight bit engine is used; we check that `inputenc` or `inputenx` has been loaded with the right option, otherwise we define `\newunicodechar` to swallow its arguments, after a warning. We need to first check for `inputenx` because it loads `inputenc`.

```

45 \def\nuc@stop{\PackageWarningNoLine{newunicodechar}
46 {This package won't work without loading\MessageBreak
47 'inputenc' or 'inputenx' with the 'utf8' option}%
48 \let\newunicodechar\@gobbletwo\endinput}
49 \@ifpackageloaded{inputenx}
50 {\def\nuc@tempa{inputenx}}
51 {\@ifpackageloaded{inputenc}{\def\nuc@tempa{inputenc}}{\nuc@stop}}
52 \@ifpackagewith{\nuc@tempa}{utf8}{\nuc@stop}
53 \@ifpackagewith{\nuc@tempa}{utf8x}{\nuc@stop}{\nuc@stop}

```

`\newunicodechar` The main macro. We set the temporary switch to false and put in `\nuc@tempa` the first argument, but detokenized, since it consists of active characters. We

check that it's not empty and access to its first token that we put into `\@tempb`. Then we call `\nuc@check` and if it sets the temporary switch to true, we execute the definition, since the first argument is a valid UTF-8 character, but we issue a warning if it already has a meaning. All we need to do is to define the macro `\csname u8:\nuc@tempa\endcsname`, because that's how `inputenc` acts. The part between `\ifdefined` and the corresponding `\fi` will be executed only with the `verbose` option.

```

54 \def\newunicodechar#1#2{%
55   \@tempswafalse
56   \edef\nuc@tempa{\detokenize{#1}}%
57   \if\relax\nuc@tempa\relax
58     \nuc@emptyargerr
59   \else
60     \edef\@tempb{\expandafter\@car\nuc@tempa\@nil}%
61     \nuc@check
62     \if@tempswa
63       \ifundefined{u8:\nuc@tempa}{%
64         {\PackageWarning{newunicodechar}
65           {Redefining Unicode character\ifdefined\nuc@verbose;
66             it meant\MessageBreak
67             ***\space\space\nuc@meaning\space\space***\MessageBreak
68             before your redefinition\fi}}%
69         \@namedef{u8:\nuc@tempa}{#2}%
70       \fi
71     \fi
72 }

```

The first helper macro computes the number of bytes in the first argument, though it appears to the user as a single character. The third is used for the `verbose` option.

```

73 \def\nuc@getlength#1{%
74   \ifx#1\@nil
75     \expandafter\relax
76   \else
77     +1\expandafter\nuc@getlength
78   \fi}
79 \ifdefined\nuc@verbose
80   \def\nuc@meaning{\expandafter\expandafter\expandafter
81     \strip@prefix\expandafter\meaning\csname u8:\nuc@tempa\endcsname}
82 \fi

```

We select the action based on the length of the first argument; in case only one byte appears, the user is trying to define an ASCII character, which is not allowed; if the input is two, three, or four bytes long, we check whether the first byte has the correct form: its binary form must be, respectively, `110xxxxx`, `1110xxxx`, or `11110xxx`, so not less than 192, 224, or 240. The macro `\nuc@ch@ck` does precisely this, issuing an error message if the argument is invalid or else setting the temporary switch to true.

```

83 \def\nuc@check{%

```

```

84 \ifcase\numexpr0\expandafter\nuc@getlength\nuc@tempa\@nil
85   \or %0
86   \nuc@onebyteerr\or %1
87   \nuc@ch@ck{192}\or %2
88   \nuc@ch@ck{224}\or %3
89   \nuc@ch@ck{240}\else %4
90   \nuc@invalidargerr
91 \fi}
92 \def\nuc@ch@ck#1{%
93   \expandafter\ifnum\expandafter'\@tempb<#1\relax
94   \nuc@invalidargerr
95   \else
96   \@tempwatrue
97 \fi
98 }

```

Finally we disallow `\newunicodechar` outside the preamble.

```
99 \@onlypreamble\newunicodechar
```

Change History

v1.0		v1.1	
	General: Initial version 1		General: Added support for inputenx 1, 5

Index

Numbers written in *italic* refer to the page where the corresponding entry is described; numbers underlined refer to the code line of the definition; numbers in *roman* refer to the code lines where the entry is used.

	Symbols	<code>\nuc@atcode</code> 23, 41	<code>\nuc@onebyteerr</code> 6, 35, 86
<code>\^</code> 19		<code>\nuc@ch@ck</code> .. 87–89, 92	<code>\nuc@stop</code> . . . 45, 51–53
<code>\~</code> 23, 24, 32, 41		<code>\nuc@check</code> 61, 83	<code>\nuc@verbose</code> .. 4, 65, 79
	D	<code>\nuc@emptyargerr</code> 11, 27, 58	P
<code>\detokenize</code> .. 26, 29, 56		<code>\nuc@getlength</code> 73, 77, 84	<code>\protected</code> 33
	N	<code>\nuc@invalidargerr</code> 14, 38, 90, 94	V
<code>\newunicodechar</code> 3, 13, 16, <u>18</u> , 48, <u>54</u> , 99		<code>\nuc@meaning</code> 67, 80	<code>verbose</code> option 2