

# Greek Unicode with 8-bit TeX and *inputenc*

Günter Milde

July 11, 2019

The definitions in `lgrenc.dfu` provide UTF-8 support for the Greek script based on *inputenc* and the *LaTeX internal character representation* macros (LICRs) defined in the *greek-fontenc* package.

## 1 Requirements

The *inputenc* standard package enables the use of non-ASCII characters with 8-bit TeX. However, it misses definitions for Greek characters. The *greek-inputenc* package extends *inputenc* to allow the use of Greek literals in the document source.

As with all *inputenc* definitions, this only works if the active font encoding supports the characters. For the Greek script, this is usually the non-standard *LGR* font encoding set up by *greek-fontenc*.

## 2 Usage

There are several alternatives to activate Greek Unicode input for 8-bit TeX<sup>1</sup> (see also the source document `greek-utf8.tex`):

- Define the LGR font encoding and the UTF8 input encoding (the order does not matter), e.g.,

```
\usepackage[T1,LGR]{fontenc}
\usepackage[utf8]{inputenc}
```

Ensure that LGR is the active font encoding whenever a Greek character is used in the text (see below).

- For text in the Greek language, it is recommended to use the *Babel* package with the Greek language definitions in *babel-greek*. Babel sets the font encoding automatically to LGR and Greek Unicode characters work as expected. Write in the preamble, e.g.,

---

<sup>1</sup> The XeTeX and LuaTeX engines use utf8 as native input encoding. They do not require (and, except in 8-bit compatibility mode, do not work with) the *inputenc* and *greek-inputenc* packages.

```

\usepackage[utf8]{inputenc}
\usepackage[LGR,T1]{fontenc}
\usepackage[english,greek,german]{babel}

```

and use `\foreignlanguage` or `\selectlanguage` to set the text language to Greek (see the [babel-greek](#) documentation for detailed examples).

Τί φήεις; Ἴδὼν ἐνθ' ἔδε παῖδ' ἔλευθέραν τὰς πλησίον Νύμφας στεφανοῦσαν,  
Σώστρατε, ἔρωϊν ἀπῆλθεσ εὐθύς;

- In combination with the *textalpha* package from [greek-fontenc](#), Greek Unicode characters can be used in text with any font encoding – just like the symbols provided by the “textcomp” package (i.e. with some limitations described in [textalpha-doc](#)). With the preamble lines

```

\usepackage[utf8]{inputenc}
\usepackage{textalpha}

```

it is straightforward to write about  $\pi$ -mesons,  $\gamma$ -radiation, or a 50 k $\Omega$  resistor.

- In combination with the *alphabet* package (also from [greek-fontenc](#)), Greek Unicode literals can also be used in math mode:

```

\usepackage[utf8]{inputenc}
\usepackage{alphabet}

```

$$\tan \beta = \frac{\sin \beta}{\cos \beta}.$$

### 3 Warning: unsafe ASCII input

LGR is no “standard font encoding”. Latin characters and some other ASCII symbols are mapped to Greek equivalents if LGR is the active font encoding. (See [usage.pdf](#) for a description of this Latin-Greek transliteration.)

This means you need an explicit language and/or font-encoding switch for Latin words and abbreviations in Greek text, e.g., not « $\eta$ ία αντίσταση 750- $\kappa\Omega$ » but « $\eta$ ία αντίσταση 750-k $\Omega$ »

Special care is also required with the question mark characters:

- The Unicode standard says character 003B SEMICOLON and not 037E GREEK QUESTION MARK, is the preferred character for a ‘Greek question mark’ (erotimatiko),
- The LGR font encoding maps a SEMICOLON to a middle dot (ano teleia), while the Latin question mark “?” is mapped to the erotimatiko.

As a result, only the deprecated character 037E GREEK QUESTION MARK works with both, Xe/LuaTeX and 8-bit TeX. Compare the source `greek-utf8.tex` and the PDF output:

Latin (T1)	Greek (LGR)	question mark character
Tí φήις;	Tí φήις;	037E GREEK QUESTION MARK
Tí φήις;	Tí φήις·	003B SEMICOLON
Tí φήις?	Tí φήις;	003F QUESTION MARK

## 4 Supported Characters

Unicode definitions exist for all non-ASCII characters that can be rendered with an LGR-encoded font.

### 4.1 Greek and Coptic

	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F
370	*	*	*	*	'		*	*				*	*	*		
380					'	´	À	·	È	Η	Í		Ό		Υ	Ω
390	ı	Α	Β	Γ	Δ	Ε	Ζ	Η	Θ	Ι	Κ	Λ	Μ	Ν	Ξ	Ο
3A0	Π	Ρ		Σ	Τ	Υ	Φ	Χ	Ψ	Ω	İ	ÿ	ά	έ	ή	ί
3B0	ύ	α	β	γ	δ	ε	ζ	η	θ	ι	κ	λ	μ	ν	ξ	ο
3C0	π	ρ	ς	σ	τ	υ	φ	χ	ψ	ω	ı	ü	ó	ú	ώ	
3D0	*	*	*	*	*	*	*	*	ϕ	ϙ	ϒ	ϛ	Ϝ	ϝ	*	ϟ
3E0	λ	λ	*	*	*	*	*	*	*	*	*	*	*	*	*	*
3F0	*	*	*	*	*	*	*	*	*	*	*	*	*	*	*	*

legend: \* glyph missing in LGR, [space] Unicode point not defined

## 4.2 Greek Extended

	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F
1F00	ά	á	â	ã	ä	å	ǎ	Ǻ	À	Á	Â	Ã	Ä	Å	Ά	Α
1F10	è	é	ê	ë	ě	Ě			È	É	Ê	Ë	Ě	Ě	Ἐ	Ἐ
1F20	ĥ	ĥ	ĥ	ĥ	ĥ	ĥ	ĥ	ĥ	Ḥ	Ḥ	Ḥ	Ḥ	Ḥ	Ḥ	Ἠ	Ἠ
1F30	ì	í	î	ï	ĭ	ĭ	ĭ	ĭ	Ì	Í	Î	Ï	Ĭ	Ĭ	Ἠ	Ἠ
1F40	ò	ó	ô	õ	ö	ő			Ò	Ó	Ô	Õ	Ö	Ö		
1F50	ù	ú	û	ü	ÿ	ÿ	ÿ	ÿ	Υ		Υ		Υ			Υ
1F60	ώ	ώ	ώ	ώ	ώ	ώ	ώ	ώ	Ω	Ω	Ω	Ω	Ω	Ω	Ω	Ω
1F70	à	á	â	ã	ä	å	ǎ	Ǻ	À	Á	Â	Ã	Ä	Å		
1F80	ą	ą	ą	ą	ą	ą	ą	ą	Ą	Ą	Ą	Ą	Ą	Ą	Ą	Ą
1F90	ĥ	ĥ	ĥ	ĥ	ĥ	ĥ	ĥ	ĥ	Ḥ	Ḥ	Ḥ	Ḥ	Ḥ	Ḥ	Ḥ	Ḥ
1FA0	ḥ	ḥ	ḥ	ḥ	ḥ	ḥ	ḥ	ḥ	Ḥ	Ḥ	Ḥ	Ḥ	Ḥ	Ḥ	Ḥ	Ḥ
1FB0	ǎ	ā	ą	ą	ą		ǎ	ą	Ā	Ā	Ā	Ā	Ā			
1FC0	˜	˜	ĥ	ĥ	ĥ		ĥ	ĥ	Ē	Ē	Ē	Ē	Ē			
1FD0	ı	ı	ı	ı			ı	ı	İ	İ	İ	İ				
1FE0	ÿ	ÿ	ÿ	ÿ	ÿ	ÿ	ÿ	ÿ	Ÿ	Ÿ	Ÿ	Ÿ	Ÿ	Ÿ	Ÿ	Ÿ
1FF0			ώ	ω	ώ		ώ	ω	Ω	Ω	Ω	Ω	Ω			

## 4.3 Other Unicode Blocks

**Latin-1 Supplement** : “ « ’ · »

**IPA Extensions** : ə LATIN SMALL LETTER SCHWA

**Spacing Modifier Letters** : ˘α (here followed by letter alpha)

**General Punctuation** : — ‘ ’ %<sub>00</sub> ZWNJ (zero width no joiner, prevents kerning and ligatures, e.g. ΑΥ vs. ΑΥ and ’α vs. ά)

**Currency Symbols** : €

**Letter-like Symbols** : Ω

**Ancient Greek Numbers** : ͵Ͷ ͵ͷ ͵͸ ͵͹

## 5 Test up/downcasing

Capital Greek letters have diacritics (except the dialytika) to the left (instead of above) and drop them in uppercase, e.g.  $\mu\acute{\alpha}\iota\sigma\tau\rho\omicron\varsigma \mapsto \text{MA}\ddot{\text{I}}\text{ΣΤΡΟΣ}$ .

Tonos and dasia on the first vowel of a diphthong ( $\acute{\alpha}\iota$ ,  $\acute{\alpha}\upsilon$ ,  $\acute{\epsilon}\iota$ ) imply a *hiatus*. A dialytika must be placed on the second vowel if they are dropped ( $\text{A}\ddot{\text{I}}$ ,  $\text{A}\ddot{\text{Y}}$ ,  $\text{E}\ddot{\text{I}}$ ).

The auto-hiatus feature in `lgrxenc.def` works with the Latin transcription and with character-macros ( $\text{A}\ddot{\text{I}}$ ,  $\text{A}\ddot{\text{Y}}$ ,  $\text{E}\ddot{\text{I}}$ ) and also if the first character is wrapped in `\ensuregreek` (as done by the `lgrxenc.dfu` definition for accented characters) or a literal Unicode character ( $\text{A}\ddot{\text{I}}$ ,  $\text{A}\ddot{\text{Y}}$ ,  $\text{A}\ddot{\text{I}}$ ) but not if the second character of the diphthong is a Unicode literal ( $\text{AI}$ ,  $\text{AY}$ ,  $\text{EI}$ ).

Therefore, the diaeresis is missing in the following examples:  $\acute{\alpha}\upsilon\lambda\omicron\varsigma \mapsto \text{AΥΛΟΣ}$ ,  $\acute{\alpha}\upsilon\lambda\omicron\varsigma \mapsto \text{AΥΛΟΣ}$ ,  $\mu\acute{\alpha}\iota\nu\alpha \mapsto \text{MAINA}$ ,  $\kappa\acute{\epsilon}\iota\kappa \mapsto \text{KEIK}$ ,  $\acute{\alpha}\upsilon\pi\nu\acute{\iota}\alpha \mapsto \text{AΥΠΝΙΑ}$ .

Fixing this shortcoming requires knowledge of what `\LGR@ifnextchar` “sees” when the next character is an upcased Unicode literal.

As an ugly workaround, use `\textiota` resp. `\textupsilon` for the character that should get the diaeresis:  $\acute{\alpha}\upsilon\pi\nu\acute{\iota}\alpha \mapsto \text{A}\ddot{\text{Y}}\text{ΠΝΙΑ}$ .

The following subsections test `MakeUppercase` and `MakeLowercase` with all characters defined in `lgrxenc.dfu`:

### 5.1 Greek and Coptic

Characters of the Greek and Coptic Unicode Block:

```
' ,; ' "A·EHTO'Y'ΩιABΓΔEZHZHΘIKAMNEOΠPΣT'YΦXΨΩİŸƆTFλ  
άέήίύαβγδεζηθικλμνξοπρςστυφχψωϊϋόύφϜϝϞϟ
```

`MakeUppercase`:

```
' ,; "A·EHIOTYŌİABΓΔEZHZHΘIKAMNEOΠPΣT'YΦXΨΩİŸƆTFλ  
AEHIŸABΓΔEZHZHΘIKAMNEOΠPΣΣT'YΦXΨΩİŸOYΩƆTFλ
```

Letters and ypogegrammeni upcased, tonos dropped, dialytika kept.

There is no capital Koppa in LGR, therefore  $\text{ϝ}$  is left unchanged with `MakeUppercase`.

`MakeLowercase`:

```
' ,; ' "ά·έήίούάαβγδεζηθικλμνξοπρστυφχψωϊϋόύφϜϝϞϟ  
άέήίύαβγδεζηθικλμνξοπρςστυφχψωϊϋόύφϜϝϞϟ
```

The lowercase of  $\Sigma$  is the «auto-sigma» (`\textautosigma`):  $\Sigma\Sigma \mapsto \sigma\varsigma$ . Add a `ZWNJ` or use the `\noboundary` macro to prevent conversion to final sigma:  $\sigma\sigma$ . The lowercase of GREEK LETTER STIGMA  $\text{Ϟ}$  is  $\text{ϟ}$ .

## 5.2 Greek extended

MakeUppercase:

Α Α Α Α Α Α Α Α Α Α Α Α Α Α Α Α  
 Ε Ε Ε Ε Ε Ε Ε Ε Ε Ε Ε Ε  
 Η Η Η Η Η Η Η Η Η Η Η Η Η Η Η Η  
 Ι Ι Ι Ι Ι Ι Ι Ι Ι Ι Ι Ι Ι Ι Ι Ι  
 Ο Ο Ο Ο Ο Ο Ο Ο Ο Ο Ο Ο  
 Υ Υ Υ Υ Υ Υ Υ Υ Υ Υ Υ Υ Υ Υ Υ Υ  
 Ω Ω Ω Ω Ω Ω Ω Ω Ω Ω Ω Ω Ω Ω Ω Ω  
 Α Α Ε Ε Η Η Ι Ι Ο Ο Υ Υ Ω Ω  
 Α<sub>1</sub> Α<sub>1</sub> Α<sub>1</sub> Α<sub>1</sub> Α<sub>1</sub> Α<sub>1</sub> Α<sub>1</sub> Α<sub>1</sub> Α<sub>1</sub> Α<sub>1</sub> Α<sub>1</sub> Α<sub>1</sub> Α<sub>1</sub> Α<sub>1</sub> Α<sub>1</sub> Α<sub>1</sub>  
 Η<sub>1</sub> Η<sub>1</sub> Η<sub>1</sub> Η<sub>1</sub> Η<sub>1</sub> Η<sub>1</sub> Η<sub>1</sub> Η<sub>1</sub> Η<sub>1</sub> Η<sub>1</sub> Η<sub>1</sub> Η<sub>1</sub> Η<sub>1</sub> Η<sub>1</sub> Η<sub>1</sub> Η<sub>1</sub>  
 Ω<sub>1</sub> Ω<sub>1</sub> Ω<sub>1</sub> Ω<sub>1</sub> Ω<sub>1</sub> Ω<sub>1</sub> Ω<sub>1</sub> Ω<sub>1</sub> Ω<sub>1</sub> Ω<sub>1</sub> Ω<sub>1</sub> Ω<sub>1</sub> Ω<sub>1</sub> Ω<sub>1</sub> Ω<sub>1</sub> Ω<sub>1</sub>  
 Ά Ά Α<sub>1</sub> Α<sub>1</sub> Α<sub>1</sub> Α Α<sub>1</sub> Ά Ά Α Α Α<sub>1</sub> ,  
 ¨ Η<sub>1</sub> Η<sub>1</sub> Η<sub>1</sub> Η Η<sub>1</sub> Ε Ε Η Η Η<sub>1</sub>  
 Ϊ Ϊ Ϊ Ϊ Ϊ Ϊ Ϊ Ϊ Ϊ Ϊ Ϊ Ϊ Ϊ Ϊ Ϊ Ϊ Ϊ Ϊ  
 ο ο ο ο Ρ Ρ Υ ο ο ο ο Υ Υ Ρ ¨ ¨  
 Ω<sub>1</sub> Ω<sub>1</sub> Ω<sub>1</sub> Ω Ω<sub>1</sub> Ο Ο Ω Ω Ω<sub>1</sub>

MakeLowercase:

α α ᾀ ᾀ ᾀ ᾀ ᾀ ᾀ ᾀ ᾀ ᾀ ᾀ ᾀ ᾀ ᾀ ᾀ  
 ε ε ε ε ε ε ε ε ε ε ε ε  
 η η η η η η η η η η η η η η η η  
 ι ι ι ι ι ι ι ι ι ι ι ι ι ι ι ι  
 ο ο ο ο ο ο ο ο ο ο ο ο ο ο ο ο  
 υ υ υ υ υ υ υ υ υ υ υ υ υ υ υ υ  
 ω ω ω ω ω ω ω ω ω ω ω ω ω ω ω ω  
 α α ε ε η η ι ι ο ο υ υ ω ω  
 ᾀ ᾀ ᾀ ᾀ ᾀ ᾀ ᾀ ᾀ ᾀ ᾀ ᾀ ᾀ ᾀ ᾀ ᾀ ᾀ  
 η η η η η η η η η η η η η η η η  
 ω ω ω ω ω ω ω ω ω ω ω ω ω ω ω ω  
 ᾱ ᾱ ᾱ ᾱ ᾱ ᾱ ᾱ ᾱ ᾱ ᾱ ᾱ ᾱ ᾱ ᾱ ᾱ  
 ~ η η η η η ε ε η η ~ ~  
 Ϊ Ϊ Ϊ Ϊ Ϊ Ϊ Ϊ Ϊ Ϊ Ϊ Ϊ Ϊ Ϊ Ϊ Ϊ Ϊ Ϊ Ϊ  
 ϰ ϰ ϰ ϰ ϰ ϰ ϰ ϰ ϰ ϰ ϰ ϰ ϰ ϰ ϰ ϰ ϰ  
 ω ω ω ω ω ω ω ω ω ω ω ω ω ω ω ω

## 5.3 Other Unicode Blocks

MakeUppercase does not change non-letter symbols and the letter shwa:

“ « - ‘ . » ə ˘ A — ‘ ’ ‰ A γ € ☒ ☒ ☒ ☒

MakeLowercase does not change non-letter symbols, too:

“ « - ‘ . » ə ˘ α — ‘ ’ ‰ α υ € ☒ ☒ ☒ ☒

