

Greek Unicode with 8-bit TeX and *inputenc*

Günter Milde

August 7, 2015

Abstract

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 TeX¹ (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.,

```
\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).

¹ 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.

Τί φής; Ἴδὼν ἐνθῆδε παῖδ' ἐλευθέραν τὰς πλησίον Νύμφας
στεφανοῦσαν, Σώστρατε, ἔρῶν 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 π-mesons, γ-radiation, or a 50 kΩ 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 «ἡία αντίσταση 750-κΩ» but «ἡία αντίσταση 750-kΩ»

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
Τί φής;	Τί φής;	037E GREEK QUESTION MARK
Τί φής;	Τί φής	003B SEMICOLON
Τί φής?	Τί φής;	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					'	´	ˆA	ˆE	ˆH	ˆI		ˆO		ˆΥ	ˆΩ	
390	ι	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	Ϊ	Ϋ	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	ϊ	ϋ	ό	ύ	ώ	
3D0	*	*	*	*	*	*	*	*	Ϟ	ϟ	Ϡ	ϡ	Ϣ	ϣ	*	ϥ
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	ά	ά̂	ά̃	ά̄	ά̅	ά̆	ά̇	Ά	Ά	Ά̂	Ά̃	Ά̄	Ά̅	Ά̆	Ά̇
1F10	έ	έ	έ̂	έ̃	έ̄	έ̅	έ̆	έ̇	Έ	Έ	Έ̂	Έ̃	Έ̄	Έ̅	Έ̆	Έ̇
1F20	ή	ή	ή̂	ή̃	ή̄	ή̅	ή̆	ή̇	Ή	Ή	Ή̂	Ή̃	Ή̄	Ή̅	Ή̆	Ή̇
1F30	ί	ί	ί̂	ί̃	ί̄	ί̅	ί̆	ί̇	Ϊ	Ϊ	Ϊ̂	Ϊ̃	Ϊ̄	Ϊ̅	Ϊ̆	Ϊ̇
1F40	ό	ό	ό̂	ό̃	ό̄	ό̅	ό̆	ό̇	Ό	Ό	Ό̂	Ό̃	Ό̄	Ό̅	Ό̆	Ό̇
1F50	ύ	ύ	ύ̂	ύ̃	ύ̄	ύ̅	ύ̆	ύ̇	Υ	Υ	Υ̂	Υ̃	Ῡ	Υ̅	Ῠ	Υ̇
1F60	ώ	ώ	ώ̂	ώ̃	ώ̄	ώ̅	ώ̆	ώ̇	Ω	Ω	Ω̂	Ω̃	Ω̄	Ω̅	Ω̆	Ω̇
1F70	ά	ά	έ	έ	ή	ή	ι	ι	ο	ο	υ	υ	ω	ω		
1F80	ά̂	ά̂	έ̂	έ̂	ή̂	ή̂	ι̂	ι̂	Α̂	Α̂	Α̂	Α̂	Α̂	Α̂	Α̂	Α̂
1F90	ή̂	ή̂	ή̂	ή̂	ή̂	ή̂	ή̂	ή̂	Η̂	Η̂	Η̂	Η̂	Η̂	Η̂	Η̂	Η̂
1FA0	ϕ̂	ϕ̂	ϕ̂	ϕ̂	ϕ̂	ϕ̂	ϕ̂	ϕ̂	Ω̂	Ω̂	Ω̂	Ω̂	Ω̂	Ω̂	Ω̂	Ω̂
1FB0	ά̃	α̃	ά̃	α̃	ά̃	α̃	ά̃	α̃	Ά̃	Α̃	Ά̃	Α̃	Α̃	Α̃	Α̃	Α̃
1FC0	~	~	ή̃	η̃	ή̃		ή̃	ή̃	Έ̃	Έ̃	Ή̃	Ή̃	Η̃	Η̃	Η̃	Η̃
1FD0	ϊ̃	ι̃	ι̃	ι̃			ι̃	ι̃	Ϊ̃	Ϊ̃	Ϊ̃	Ϊ̃	Ϊ̃	Ϊ̃	Ϊ̃	Ϊ̃
1FE0	ϋ̃	υ̃	υ̃	υ̃	ρ̃	ρ̃	υ̃	υ̃	Υ̃	Υ̃	Υ̃	Υ̃	Υ̃	Υ̃	Υ̃	Υ̃
1FF0			ϕ̃	φ̃	ϕ̃		ω̃	ω̃	Ό̃	Ό̃	Ό̃	Ό̃	Ό̃	Ό̃	Ό̃	Ό̃

4.3 Other Unicode Blocks

Latin-1 Supplement : “ « ’ · »

IPA Extensions : ə LATIN SMALL LETTER SCHWA

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

General Punctuation : – — ‘ ’ ‰ 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 : ͵Ͷͷ͸͹

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. $\mu\acute{\alpha}\sigma\tau\rho\omicron\varsigma \mapsto \text{MA-}\acute{\iota}\text{ΣΤΡΟΣ}$.

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 ($\text{A}\acute{\iota}$, $\text{A}\acute{\upsilon}$, $\text{E}\acute{\iota}$).

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

Therefore, the diaeresis is missing in the following examples: $\acute{\alpha}\upsilon\lambda\omicron\varsigma \mapsto \text{A}\acute{\upsilon}\lambda\omicron\text{Σ}$, $\mu\acute{\alpha}\iota\nu\alpha \mapsto \text{M}\text{A}\text{I}\text{N}\text{A}$, $\kappa\acute{\epsilon}\iota\kappa \mapsto \text{K}\text{E}\text{I}\text{K}$, $\acute{\alpha}\upsilon\pi\nu\acute{\iota}\alpha \mapsto \text{A}\acute{\upsilon}\text{P}\text{I}\text{N}\text{I}\text{A}$.

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: $\acute{\alpha}\upsilon\pi\nu\acute{\iota}\alpha \mapsto \text{A}\acute{\iota}\text{P}\text{I}\text{N}\text{I}\text{A}$.

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

5.1 Greek and Coptic

Characters of the Greek and Coptic Unicode Block:

```
’; ’ “A·E’HTO’T’Ω’ΑΒΓΔΕΖΗΘΙΚΑΜΝΞΟΠΡΣΤΥΦΧΨΩΪΫϞϟϠ  
άέήίύάβγδεζηθικλμνξοπρστυφχψωϊϋόύϙϚϛ
```

MakeUppercase:

```
’; ; ’ “A·EHIOTYΩ’ΑΒΓΔΕΖΗΘΙΚΑΜΝΞΟΠΡΣΤΥΦΧΨΩΪΫϞϟϠ  
ΑΕΗΙΎΑΒΓΔΕΖΗΘΙΚΑΜΝΞΟΠΡΣΣΤΥΦΧΨΩΪΎΟΥΩϞϟϠ
```

Letters and ypogegrammeni upcased, tonos dropped, dialytika kept.

There is no capital Koppa in LGR, therefore ϣ is left unchanged with `MakeUppercase`.

MakeLowercase:

```
’; ’ “ά·έήίούώάβγδεζηθικλμνξοπρστυφχψωϊϋόύϙϚϛ  
άέήίύάβγδεζηθικλμνξοπρστυφχψωϊϋόύϙϚϛ
```

The lowercase of Σ is the «auto-sigma» (`\textautosigma`): $\Sigma \mapsto \sigma\varsigma$.

Add a `ZWNJ` or use the `\noboundary` macro to prevent conversion to final sigma: $\sigma\sigma$. The lowercase of GREEK LETTER STIGMA Ϟ is τ .

5.2 Greek extended

MakeUppercase:

```
A A A A A A A A A A A A A A A A  
E E E E E E E E E E  
H H H H H H H H H H H H H H  
I I I I I I I I I I I I I I  
O O O O O O O O O O O  
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  
A A E E H H I I O O Υ Υ Ω Ω
```

A₁ A₁ A₁ A₁ A₁ A₁ A₁ A₁ A₁ A₁ A₁ A₁ A₁ A₁ A₁ A₁ A₁
 H₁ H₁ H₁ H₁ H₁ H₁ H₁ H₁ H₁ H₁ H₁ H₁ H₁ H₁ H₁ H₁ H₁
 Ω₁ Ω₁ Ω₁ Ω₁ Ω₁ Ω₁ Ω₁ Ω₁ Ω₁ Ω₁ Ω₁ Ω₁ Ω₁ Ω₁ Ω₁ Ω₁
 Ā Ā A₁ A₁ A₁ A₁ A₁ Ā Ā A₁ A₁ Ā
 “ H₁ H₁ H₁ H₁ E E H H H₁
 Ī Ī Ī Ī Ī Ī Ī Ī Ī
 Ŷ Ŷ Ŷ Ŷ P P Ŷ Ŷ Ŷ Ŷ P P …
 Ω₁ Ω₁ Ω₁ Ω₁ Ω₁ O O Ω Ω Ω₁

MakeLowercase:
 ǎ ǎ ǎ ǎ ǎ ǎ ǎ ǎ ǎ ǎ ǎ ǎ
 è è è è è è è è è è
 ħ ħ ħ ħ ħ ħ ħ ħ ħ ħ ħ ħ
 ï ï ï ï ï ï ï ï ï ï ï ï
 ó ó ò ò ò ò ó ó ò ò
 ù ù ù ù ù ù ù ù ù ù ù ù
 ǎ ǎ è é ħ ħ ì ì ò ó ù ú ǎ ǎ
 ǎ ǎ ǎ ǎ ǎ ǎ ǎ ǎ ǎ ǎ ǎ ǎ
 ħ ħ ħ ħ ħ ħ ħ ħ ħ ħ ħ ħ
 ǎ ǎ ǎ ǎ ǎ ǎ ǎ ǎ ǎ ǎ ǎ ǎ
 ǎ ǎ ǎ ǎ ǎ ǎ ǎ ǎ ǎ ǎ ǎ ǎ
 ~ ħ ħ ħ ħ è é ħ ħ ħ ħ
 ŷ ŷ ŷ ŷ ŷ ŷ ŷ ŷ ŷ ŷ ŷ ŷ
 ŷ ŷ ù ù ù ù ŷ ŷ ŷ ŷ ù ù ù ù
 ǎ ǎ ǎ ǎ ǎ ǎ ǎ ǎ ǎ ǎ ǎ ǎ

5.3 Other Unicode Blocks

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

“ « - ’ . » ə ˘ A — ‘ ’ % 0 A ʔ € ☒ ☒ ☒ ☒

MakeLowercase does not change non-letter symbols, too:

“ « - ’ . » ə ˘ α — ‘ ’ % 0 α υ € ☒ ☒ ☒ ☒

6 Test kerning/ligatures

check for kerning and unwanted ligatures:

Aǎα Aǎα Aǎα Aǎα Aǎα Aǎα Aǎα Aǎα Aǎα Aǎα Aǎα Aǎα Aǎα Aǎα Aǎα
 AʹAα AʹAα AʹAα AʹAα
 Aêα Aêα Aêα Aêα Aêα Aêα AʹEα AʹEα AʹEα AʹEα AʹEα AʹEα AʹEα AʹEα
 Aĥα Aĥα Aĥα Aĥα Aĥα Aĥα Aĥα Aĥα AʹHα AʹHα AʹHα AʹHα AʹHα AʹHα
 AʹHα AʹHα AʹHα AʹHα
 Aíα Aíα Aíα Aíα Aíα Aíα Aíα Aíα Aíα Aíα Aíα Aíα Aíα Aíα Aíα Aíα
 Aòα Aòα Aòα Aòα Aòα Aòα AʹOα AʹOα AʹOα AʹOα AʹOα AʹOα AʹOα AʹOα
 Aùα Aùα Aùα Aùα Aùα Aùα Aùα Aùα AʹUα AʹUα AʹUα AʹUα AʹUα AʹUα
 Aóα Aóα Aóα Aóα Aóα Aóα Aóα Aóα AʹΩα AʹΩα AʹΩα AʹΩα AʹΩα AʹΩα
 AʹΩα AʹΩα AʹΩα AʹΩα
 Aàα Aáα Aêα Aéα Aĥα Aĥα Aíα Aíα Aòα Aóα Aùα Aùα Aóα
 Aóα
 Aǎα Aǎα Aǎα Aǎα Aǎα Aǎα Aǎα Aǎα Aǎα Aǎα Aǎα Aǎα Aǎα Aǎα Aǎα
 AʹA₁α AʹA₁α AʹA₁α AʹA₁α

$A\grave{\eta}\alpha$ $A\acute{\eta}\alpha$ $A\tilde{\eta}\alpha$ $A\bar{\eta}\alpha$ $A\check{\eta}\alpha$ $A\dot{\eta}\alpha$ $A\ddot{\eta}\alpha$ $A\text{'}\eta_1\alpha$ $A\text{'}\eta_1\alpha$ $A^{\text{''}}\eta_1\alpha$ $A^{\text{''}}\eta_1\alpha$
 $A^{\text{''}}\eta_1\alpha$ $A^{\text{''}}\eta_1\alpha$ $A^{\text{''}}\eta_1\alpha$ $A^{\text{''}}\eta_1\alpha$
 $A\grave{\phi}\alpha$ $A\acute{\phi}\alpha$ $A\tilde{\phi}\alpha$ $A\bar{\phi}\alpha$ $A\check{\phi}\alpha$ $A\dot{\phi}\alpha$ $A\ddot{\phi}\alpha$ $A\text{'}\phi_1\alpha$ $A\text{'}\phi_1\alpha$ $A^{\text{''}}\phi_1\alpha$
 $A^{\text{''}}\phi_1\alpha$ $A^{\text{''}}\phi_1\alpha$ $A^{\text{''}}\phi_1\alpha$ $A^{\text{''}}\phi_1\alpha$ $A^{\text{''}}\phi_1\alpha$
 $A\check{\alpha}$ $A\bar{\alpha}$ $A\dot{\alpha}$ $A\alpha$ $A\acute{\alpha}$ $A\check{\alpha}$ $A\tilde{\alpha}$ $A\bar{\alpha}$ $A\text{'}\alpha$ $A\text{'}\alpha$ $AA_1\alpha$
 $A\text{'}\alpha$ $A_1\alpha$ $A\text{'}\alpha$
 $A\tilde{\alpha}$ $A^{\text{''}}\alpha$ $A\grave{\eta}\alpha$ $A\eta_1\alpha$ $A\acute{\eta}\alpha$ $A\tilde{\eta}\alpha$ $A\bar{\eta}\alpha$ $A\text{'}\text{E}\alpha$ $A\text{'}\text{E}\alpha$ $A\text{'}\text{H}\alpha$ $A\text{'}\text{H}\alpha$ $A\text{H}_1\alpha$
 $A^{\text{''}}\alpha$ $A^{\text{''}}\alpha$ $A^{\text{''}}\alpha$
 $A\text{'}\alpha$ $A\bar{\alpha}$ $A\text{'}\alpha$ $A\text{'}\alpha$ $A\text{'}\alpha$ $A\text{'}\alpha$ $A\text{'}\alpha$ $A\text{'}\alpha$ $A\text{'}\alpha$ $A\text{'}\alpha$ $A\text{'}\alpha$ $A\text{'}\alpha$ $A\text{'}\alpha$ $A\text{'}\alpha$
 $A\check{\alpha}$ $A\bar{\alpha}$ $A\text{'}\alpha$ $A\text{'}\alpha$ $A\text{'}\alpha$ $A\text{'}\alpha$ $A\text{'}\alpha$ $A\text{'}\alpha$ $A\text{'}\alpha$ $A\text{'}\alpha$ $A\text{'}\alpha$ $A\text{'}\alpha$ $A\text{'}\alpha$ $A\text{'}\alpha$
 $A\check{\rho}\alpha$ $A\dot{\rho}\alpha$ $A\ddot{\rho}\alpha$ $A\text{'}\rho_1\alpha$ $A\text{'}\rho_1\alpha$ $A^{\text{''}}\rho_1\alpha$ $A^{\text{''}}\rho_1\alpha$ $A\text{'}\Upsilon\alpha$ $A\text{'}\Upsilon\alpha$ $A\text{'}\Upsilon\alpha$ $A\text{'}\Upsilon\alpha$
 $A\text{'}\text{P}\alpha$ $A\text{'}\alpha$ $A\text{'}\alpha$ $A\text{'}\alpha$
 $A\grave{\omega}\alpha$ $A\acute{\omega}\alpha$ $A\tilde{\omega}\alpha$ $A\bar{\omega}\alpha$ $A\check{\omega}\alpha$ $A\text{'}\text{O}\alpha$ $A\text{'}\text{O}\alpha$ $A\text{'}\Omega\alpha$ $A\text{'}\Omega\alpha$ $A\Omega_1\alpha$ $A\text{'}\alpha$ $A\text{'}\alpha$