

Greek Unicode with 8-bit TeX and *inputenc*

Günter Milde

March 2, 2023

Abstract

The definitions in `1grecenc.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.,

¹The 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 φήις; Ιδών ἐνθέδε παιᾶν ἔλευθέρων τὰς πλησίου Νύμφας στεφανοῦσαν, Σώστρατε, ἐρῶν ἀπῆλθες εὐθύς;

- Greek literal characters can also be used in PDF-strings (bookmarks and ToC entries with *hyperref*). See *greek-fontenc* for a *hyperref* test and usage example.
- 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 $50\text{ k}\Omega$ resistor.²

- 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 «ηία αντίσταση 750-κΩ» but «ηία 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’ (*erotimatiiko*),

²The MIKRO SIGN and OHM SIGN characters are set up by *inputenc* with *textcomp* characters for any font encoding while GREEK CAPITAL LETTER OMEGA works only with the LGR font encoding.

- The LGR font encoding maps a SEMICOLON to a middle dot (ano teleia), while the Latin question mark “?” is mapped to the erotimatiiko.

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

With the *textalpha* package’s “keep-semicolon” option, the SEMICOLON character can be used for the erotimatiiko also with LGR encoded fonts.

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	Δ	E	Φ
370	*	*	*	*	'	*	*	*	*	*	*	*	*	*	;	
380					'	'	'	'	'	'	'	'	'	'	'	
390	τ	Α	Β	Γ	Δ	Ε	Ζ	Η	Θ	Ι	Κ	Λ	Μ	Ν	Ξ	
3A0	Π	Ρ		Σ	Τ	Υ	Φ	Χ	Ψ	Ω	Ϊ	Ӯ	ӻ	ӻ	ӻ	
3B0	ύ	α	β	γ	δ	ε	ζ	η	θ	ι	κ	λ	μ	ν	ξ	
3C0	π	ρ	ς	σ	τ	υ	φ	χ	ψ	ω	Ӯ	ӻ	ӻ	ӻ	ӻ	
3D0	*	*	*	*	*	*	*	*	*	*	Ͳ	Ͳ	F	F	*	
3E0	Δ	Ϟ	*	*	*	*	*	*	*	*	*	*	*	*	*	
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	Δ	E	Φ
1F00	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ
1F10	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ
1F20	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ
1F30	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ
1F40	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ
1F50	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ
1F60	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ
1F70	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ
1F80	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ
1F90	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ
1FA0	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ
1FB0	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ
1FC0	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ
1FD0	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ
1FE0	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ
1FF0	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ	ጀ

4.3 Other Unicode Blocks

Latin-1 Supplement : “ « ‘ ’ »

IPA Extensions : ø 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. AY vs. AY and ’α vs. á)

Currency Symbols : €

Letter-like Symbols : Ω

Ancient Greek Numbers : Ι ΙΙ ΙΙΙ ΙΙΙΙ ΙΙΙΙΙ ΙΙΙΙΙΙ

5 Test up/downcasing

Capital Greek letters have diacritics (except the dialytika, macron, and breve) to the left (instead of above) and drop them in uppercase, e.g. μαῖστρος → ΜΑΪΣΤΡΟΣ.

Tonos and dasia on the first vowel of a diphthong (άι, ῥι, εί) imply a *hiatus*. A dialytika must be placed on the second vowel if they are dropped (ΑΪ, ΑΪ̄, ΕΪ̄).

The implementation of \MakeUppercase changed in the 2022/06 LaTeX release (cf. LaTeX News 35). Since then, Greek uppercase rules are only applied if the text language is set to “greek” with Babel.

With the “old” \MakeUppercase, the auto-hiatus feature in LGR works with the Latin transliteration and LICR macros (ΑΪ, ΑΪ̄, ΕΪ̄) but not with literal Unicode characters.

It works also if only the first character of the diphthong is wrapped in \ensuregreek (as done by the lgrenc.dfu definition for accented characters)) or input as literal Unicode character (AΪ, AӮ, AӬ) but not if the second character of the diphthong is a Unicode literal (AI, AF, EI).

As an ugly workaround, use \textiota resp. \textupsilon for the character that should get the diaresis: $\grave{\alpha}\nu\tau\acute{\nu}\alpha$ \leftrightarrow ΑΓΓΙΝΙΑ.

With the “new” (post 2022/06) \MakeUppercase, the auto-hiatus feature in LGR works with literals from the “Greek and Coptic” Unicode block and some from “Greek extended”. See the following examples: ἀλος ↔ ΑΥΛΟΣ, ἀυλος ↔ ΑΤΑΛΟΣ, μαίνα ↔ MAINA, κέικ, ↔ KEIK, ἀπνία ↔ ΑΠΝΙΑ.

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:

„Α·Ε·Η·Τ·Ο·Υ·Ω·Α·Β·Γ·Δ·Ε·Ζ·Η·Ι·Κ·Λ·Μ·Ν·Ξ·Ο·Π·Ρ·Σ·Τ·Υ·Φ·Χ·Ψ·Ω·Ϊ·Γ·Ω·Γ·Φ·Λ·
ά·έ·ή·ή·ή·β·γ·δ·ε·ζ·η·θ宣·ι·λ·μ·ν·ξ宣·π·ρ·ε·σ·τ宣·ψ·ω宣·ύ宣·ό宣·ώ宣·φ·τ宣·ή宣·λ宣·

MakeUppercase:

ΑΕΗΙΤΑΒΓΔΕΖΗΘΙΚΑΛΜΝΞΟΠΡΣΤΥΦΧΨΩΪΤΩΔΦΛ
ΑΕΗΙΤΑΒΓΔΕΖΗΘΙΚΑΛΜΝΞΟΠΡΣΣΤΥΦΧΨΩΪΤΩΔΦΛ

Letters and yπογεραμμενι upcased, tonos dropped, dialytika kept.

There is no capital Koppa in LGR, therefore ' is left unchanged with Make-Uppercase.

MakeLowercase:

‘;¹ ‘ά·έ̄ή̄ίόύώ̄αβγδεζηθικλμνξοπρστυφχψω̄ϋ̄φτφλ
ά·έ̄ή̄ίύαβγδεζηθικλμνξοπρστυφγψω̄ϋ̄όύώ̄φτφλ

The lowercase of Σ 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\zeta`. The lowercase of GREEK LETTER STIGMA Υ is ξ .

5.2 Greek extended

MakeUppercase:

MakeLowercase:

5.2.1 Iota subscript vs. iota adscript

Pre-composed capital letters with *mute iota* decompose to the base letter and COMBINING GREEK YPOGEGRAMMENI (U+0345) even if they are named ... WITH [...] AND] PROSGEGRAMMENI for “historic reasons”.³

Accordingly, the “canonical” LICR for capital letters with mute iota is the base character LICR followed by \ypogegrammeni.

Compare letters followed by \prosgegrammeni and \ypogegrammeni with the pre-composed characters and with character + literal GREEK YPOGEGRAMMENI.

prosgegrammeni: $\alpha_i \alpha_i \alpha \alpha_i / A_i A_i A_i A_i / 'A_i 'A_i 'A_i 'A_i$
 MakeUppercase $A_i A_i A_i A_i / A_i A_i A_i A_i / A_i A_i A_i A_i$
 MakeLowercase $\alpha \alpha \alpha \alpha / \alpha \alpha \alpha \alpha / \dot{\alpha} \dot{\alpha} \dot{\alpha} \dot{\alpha}$
 yopogrammeni: $\alpha \alpha \alpha \alpha_i / A_i A_i A_i A_i / 'A_i 'A_i 'A_i 'A_i$
 MakeUppercase $A_i A_i A_i A_i / A_i A_i A_i A_i / A_i A_i A_i A_i$
 MakeLowercase $\alpha \alpha \alpha \alpha / \alpha \alpha \alpha \alpha / \dot{\alpha} \dot{\alpha} \dot{\alpha} \dot{\alpha}$

5.3 Other Unicode Blocks

`MakeUppercase` does not change non-letter symbols and the letter shwa (there is a capital Cyrillic schwa in T2A encoded fonts):

“《—’·》ア^A=—‘% A^€ A H X M

³cf. Nick Nicholas ‘Titlecase and Adscripts’

`MakeLowercase` does not change non-letter symbols either:

“《—·》ə ‘α — ‘ ’ %o αυ € Α Η Ι Μ

6 Test kerning/ligatures

check for kerning and unwanted ligatures: